

Архитектура надёжной In-Memory-СУБД на примере Tarantool

Владимир Перепелица



HighLoad⁺⁺
2022

Яндекс



Архитектура надёжной In-Memory-СУБД на примере Tarantool

Mons Anderson



HighLoad⁺⁺
2022

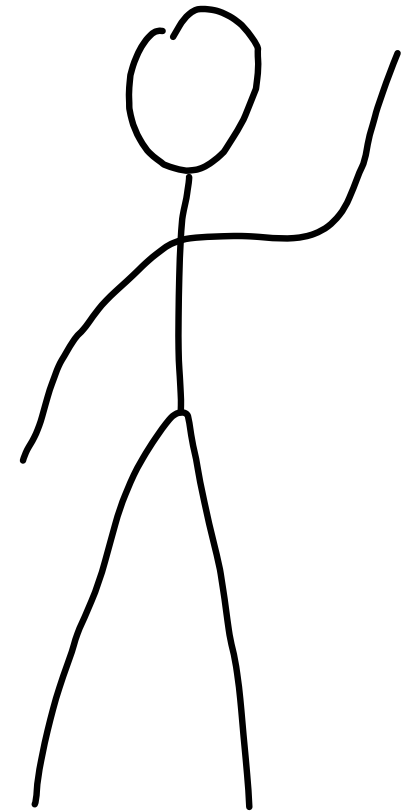
Яндекс



Кто я?

Архитектор Облака Mail.ru, VK Cloud (S3)

Архитектор и продакт-менеджер Tarantool

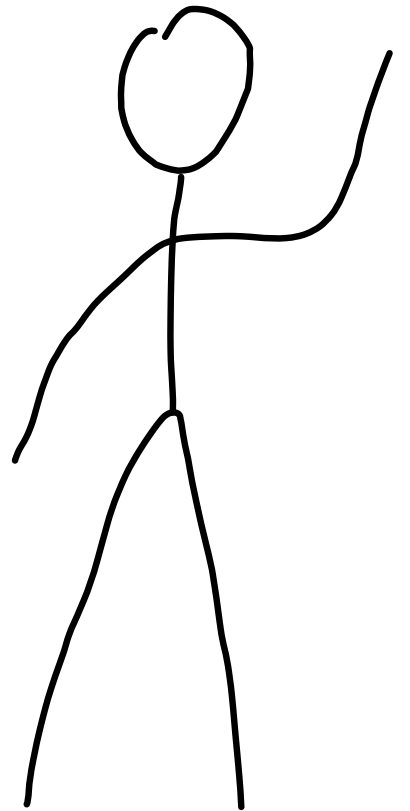


Кто я?

Архитектор Облака Mail.ru, VK Cloud (S3)

Архитектор и продакт-менеджер Tarantool

Евангелист Tarantool



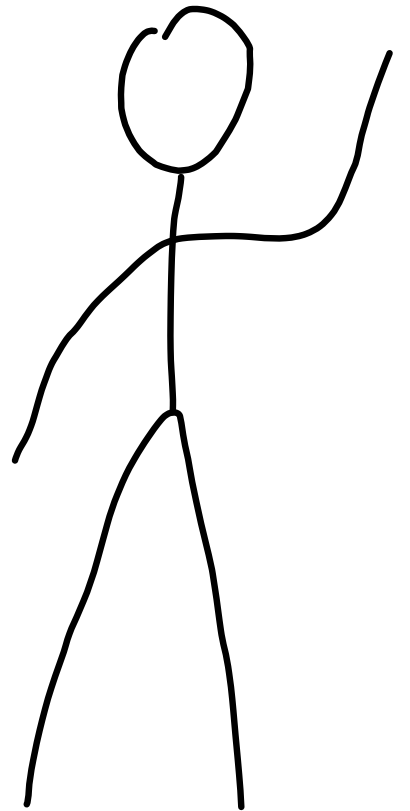
Кто я?

Архитектор Облака Mail.ru, VK Cloud (S3)

Архитектор и продакт-менеджер Tarantool

Евангелист Tarantool

Использую Tarantool более 10 лет



Кто я?

Архитектор Облака Mail.ru, VK Cloud (S3)

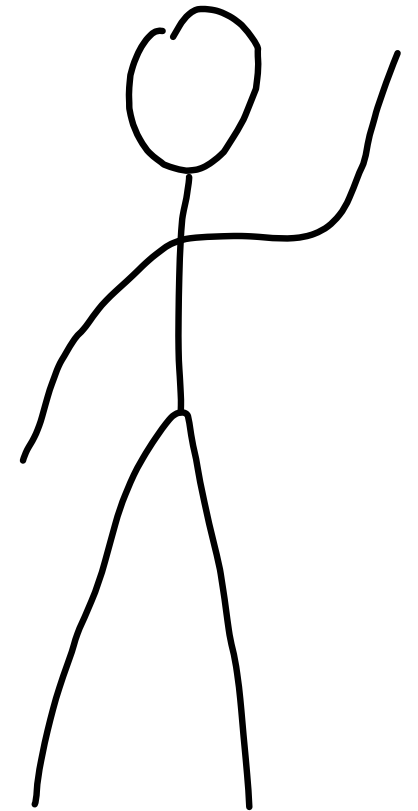
Архитектор и продакт-менеджер Tarantool

Евангелист Tarantool

Использую Tarantool более 10 лет

Как основное хранилище

Без «запасного» варианта рядом



О чём мы поговорим?

01 Почему In-Memory?

02 От кеша к БД

03 Персистентность

04 Репликация

05 MVCC

О чём мы поговорим?

01 Почему In-Memory?

02 От кеша к БД

03 Персистентность

04 Репликация

05 MVCC

Ссылка на слайды будет в конце

01

Зачем нужны In-Memory СУБД?

И почему не взять «старый добрый» SQL с кешом на тер?

Откуда взялись In-Memory БД?

- Изначально БД — дисковые
- И медленные
- Память — это кеш

Откуда взялись In-Memory БД?

- Изначально БД — дисковые
- И медленные
- Память — это кеш
- В начале 00-х производительности стало не хватать
- Появились первые решения для хранения в памяти

TimesTen: 1998 (!)

In-Memory RDBMS

Разработана в HP

Куплена Oracle в 2006

Проприетарная :(

137. ↓	132. ↓	122.	TimesTen	Relational	2.02	-0.07	-0.03
--------	--------	------	----------	------------	------	-------	-------

Memcached: 2003

Распределённая система кеширования в памяти
Для Livejournal

Простой API

Шардинг на клиенте

Автовытеснение

Кеш, а не надёжное хранение



Кто ещё?

TimesTen 1998

Memcached 2003

Hazelcast 2008

Tarantool 2008

Redis 2009

Aerospike 2012

Apache Ignite 2015



Кто ещё?

TimesTen 1998 (C)

Memcached 2003 (Perl/C)

Hazelcast 2008 (Java)

Tarantool 2008 (C)

Redis 2009 (C)

Aerospike 2012 (C)

Apache Ignite 2015 (Java)



Почему не <любимый-SQL> с кешом?

- Обычные реляционные БД ориентированы на диск
- Работа с диском не экономит размер
- Работа с диском экономит I/O(ps)
- Работа с диском требует многопоточности
- Многопоточность требует блокировок

Что можно сделать особенного?

- Монопольный доступ к памяти (1 поток)
- Оптимизация под память
- Специализированные аллокаторы
- Ориентация на OLTP (K/V)

Выводы: зачем In-Memory?

Потому что нужна производительность!

Большинство решений даёт сотни тысяч OLTP с ядра.

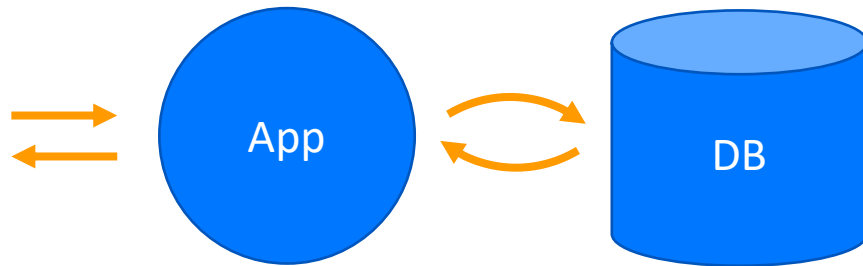
Пример:

Tarantool развивает до 1M RPS смешанных операций

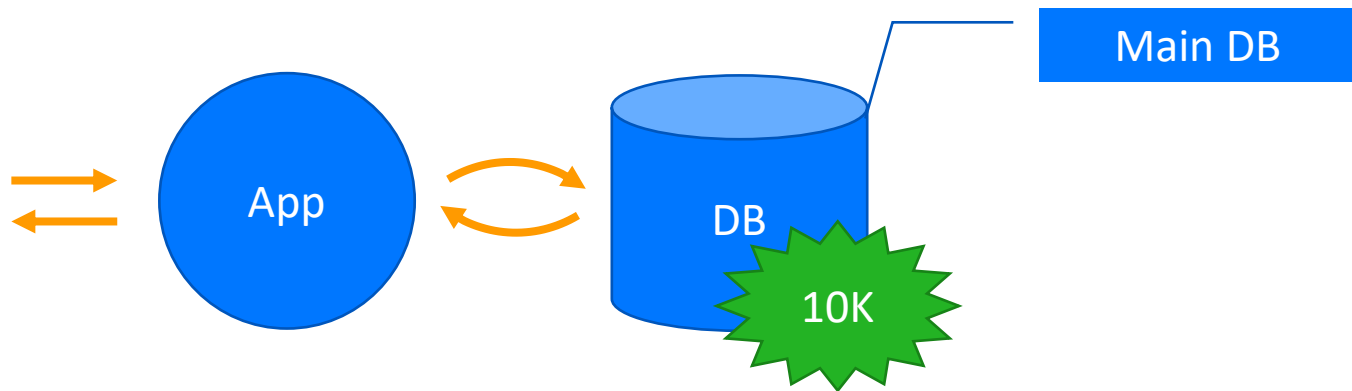
За счёт чего?

- Нет чтения диска
- Запись транзакций линейная
- Нет необходимости кэшировать
- Идеология массового обслуживания

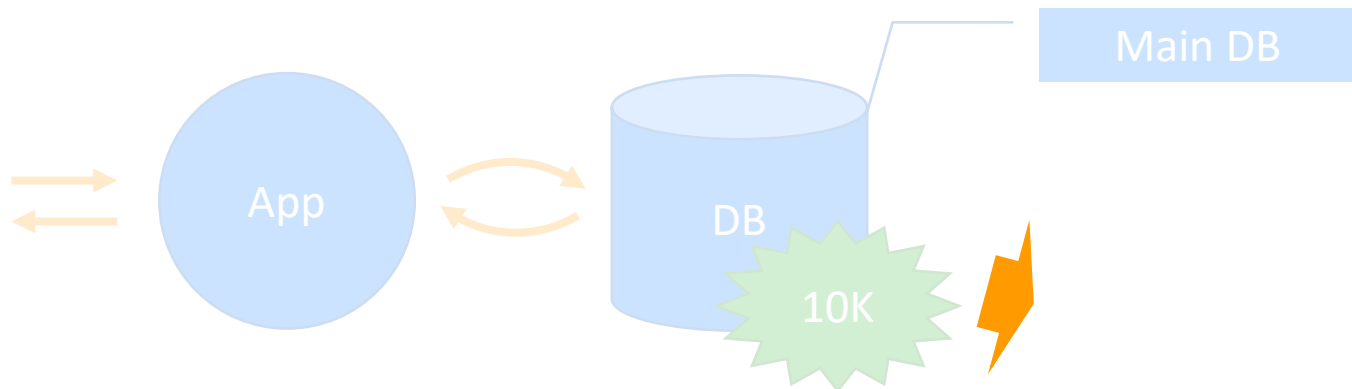
Выводы: зачем In-Memory?



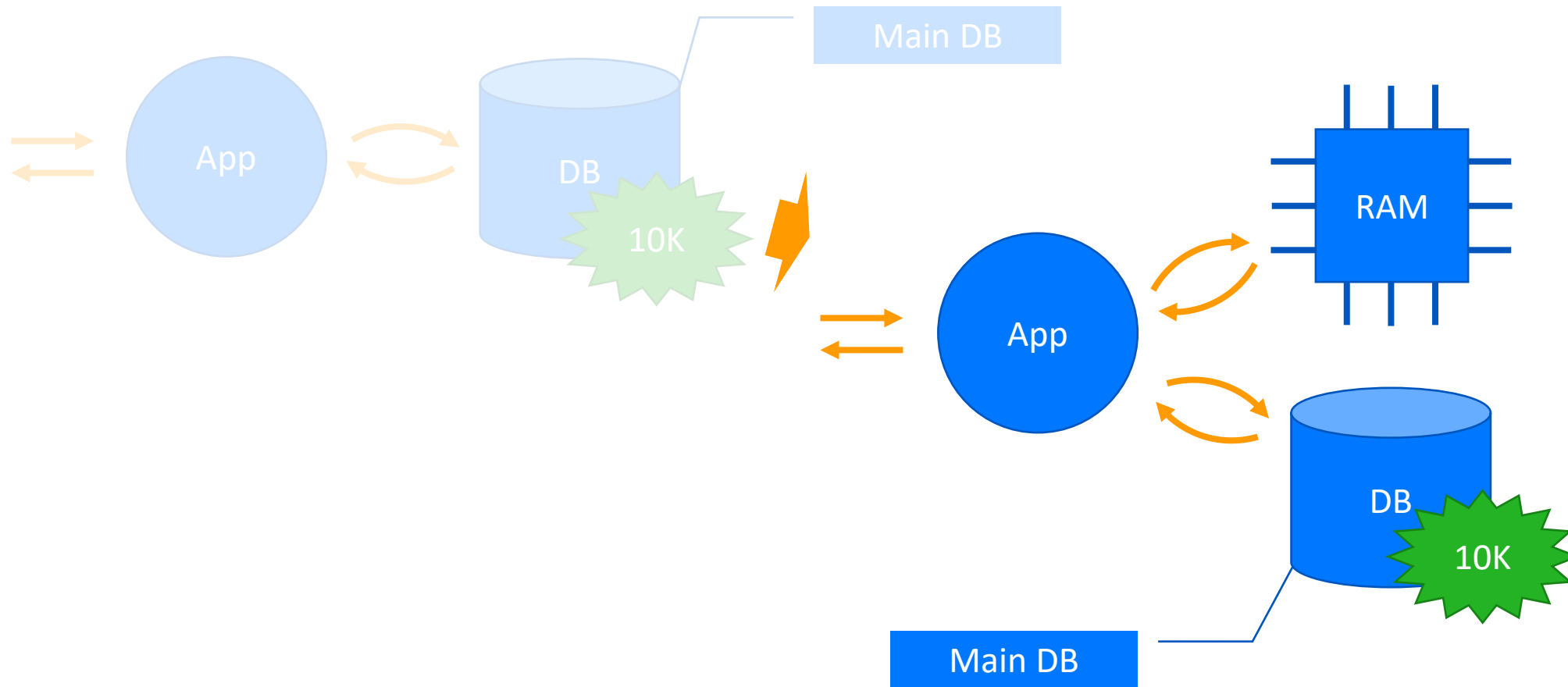
Выводы: зачем In-Memory?



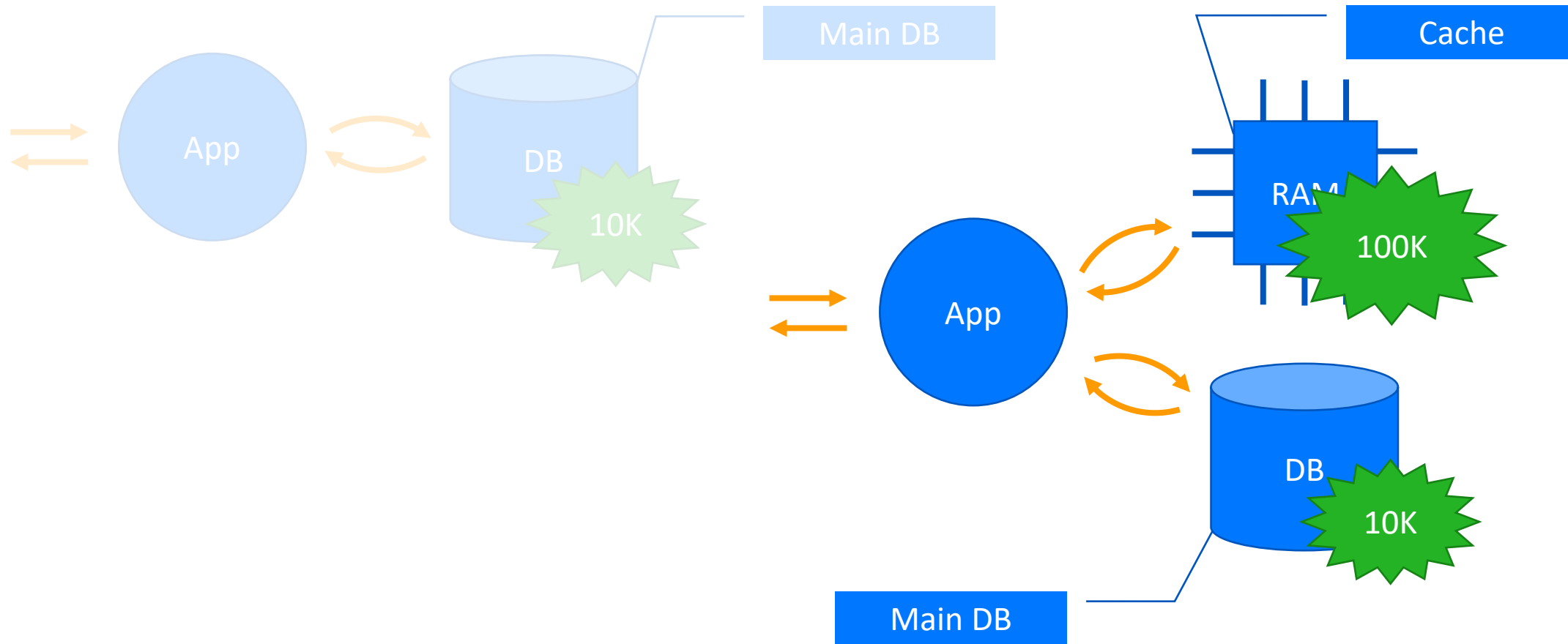
Выводы: зачем In-Memory?



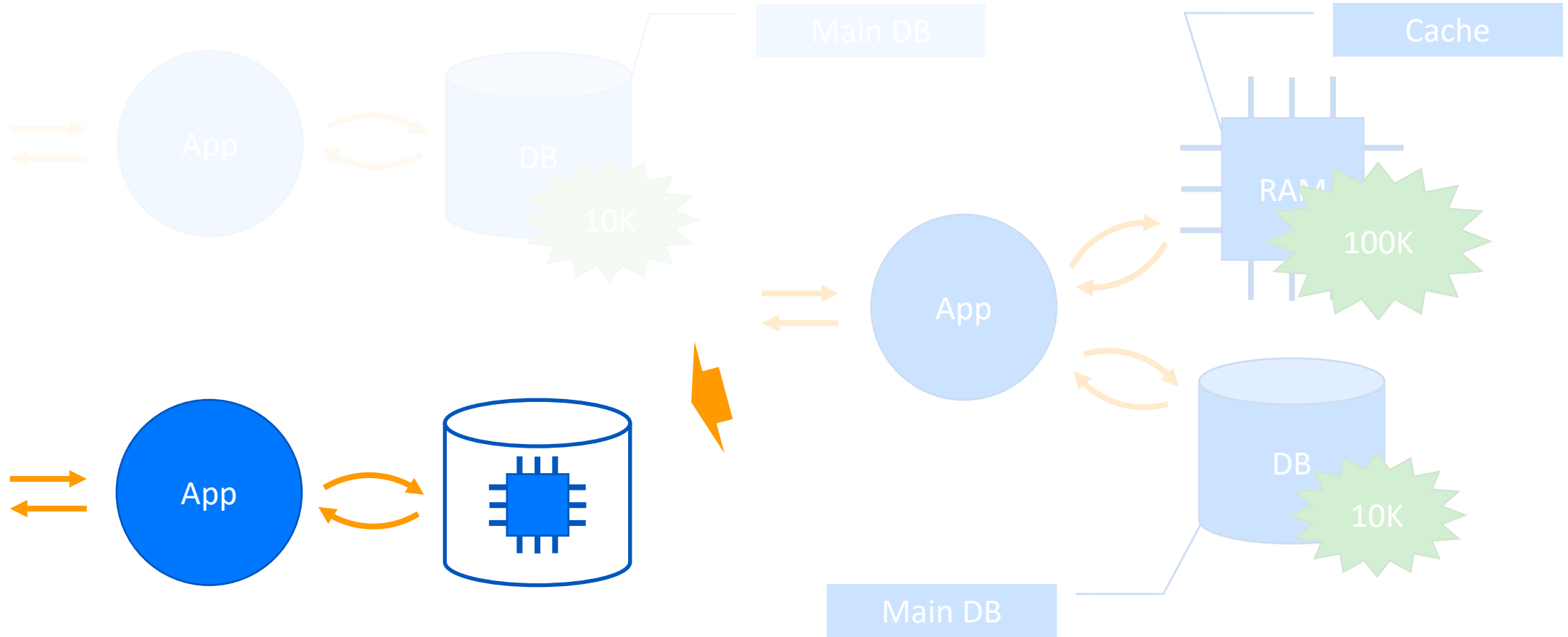
Выводы: зачем In-Memory?



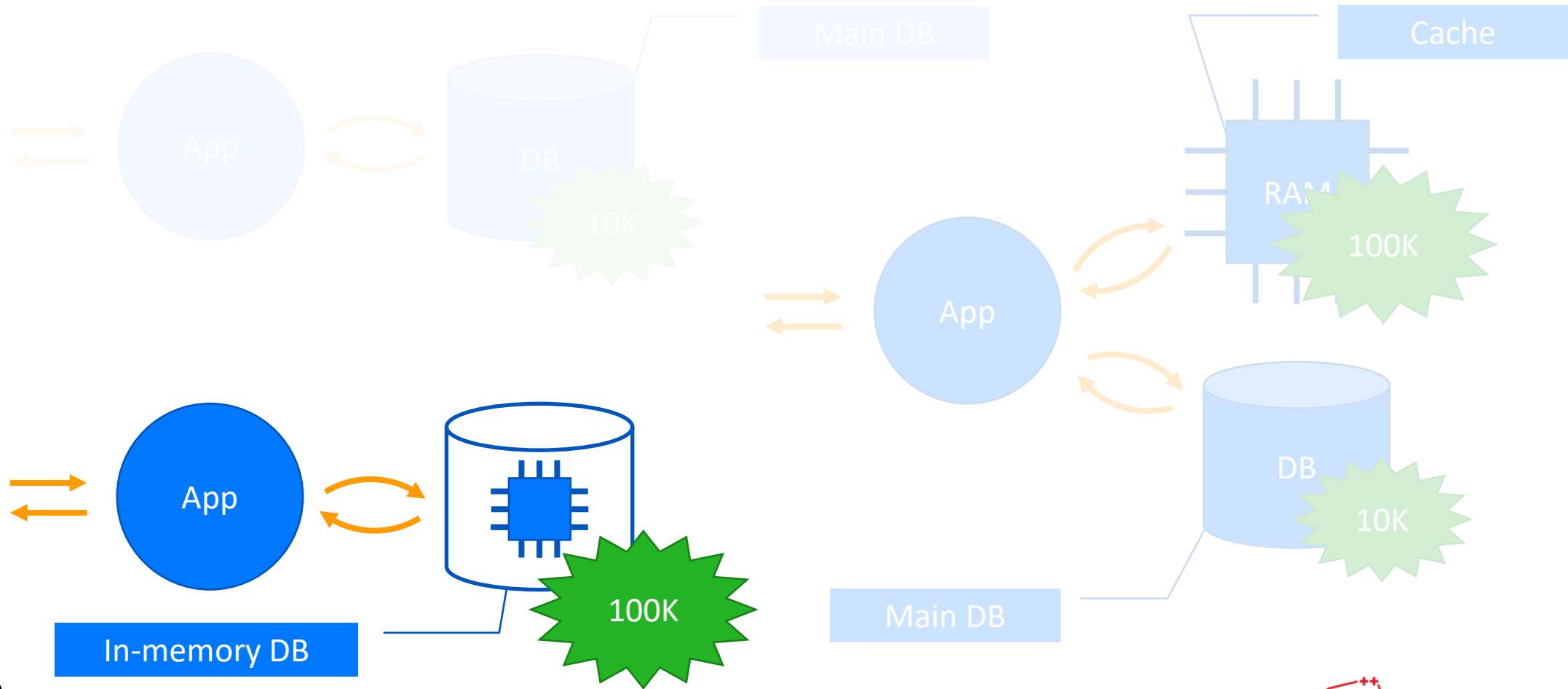
Выводы: зачем In-Memory?



Выводы: зачем In-Memory?



Выводы: зачем In-Memory?



Но какой ценой?



Did you do it?



Yes.



What did it cost?



Everything.



**Everything.
In RAM.**

02

От кеша к БД

Что отличает In-Memory БД от кеша в памяти?

История Tarantool

База данных

Кеш со вторичными индексами

Очередь

Платформа in-memory вычислений

База данных и сервер приложений

Мемкеш на стероидах

Резидентная БД

Сервер приложений Lua

Как Redis, но лучше



История Tarantool

in-memory cache

История Tarantool

in-memory *persistent* cache

История Tarantool



in-memory persistent cache

in-memory persistent cache with replication

2008

История Tarantool



in-memory persistent cache

2008

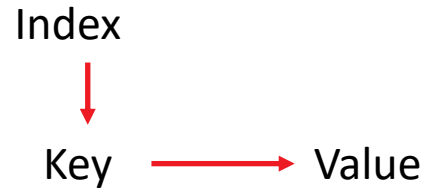
in-memory persistent cache with replication

in-memory key/value database

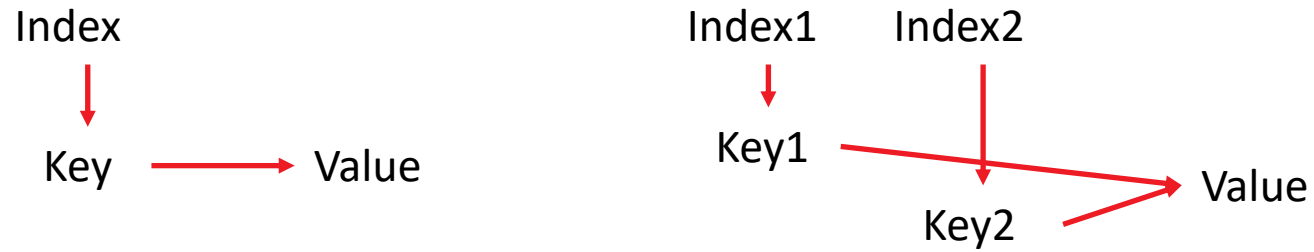
2010

in-memory multi-index database (ACID)

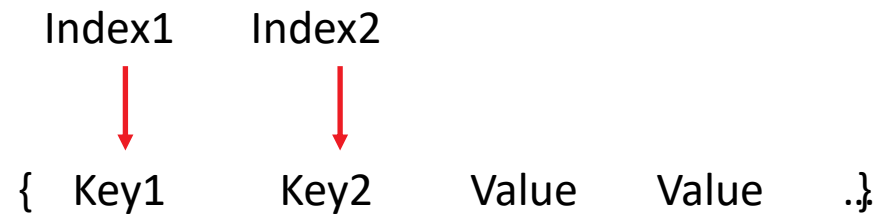
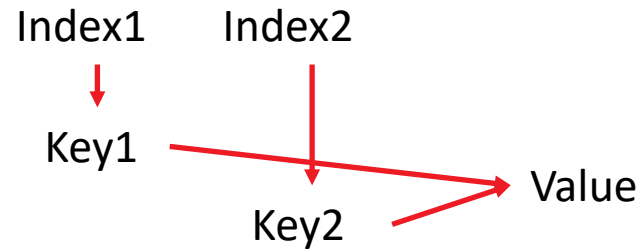
Key/value to Multi-index



Key/value to Multi-index



Key/value to Multi-index



История Tarantool



in-memory persistent cache

2008

in-memory persistent cache with replication

in-memory key/value database

2010

in-memory multi-index database (ACID)

in-memory multi-index db with lua functions

2012

История Tarantool



in-memory persistent cache

2008

in-memory persistent cache with replication

in-memory key/value database

2010

in-memory multi-index database (ACID)

in-memory multi-index db with lua functions

2012

in-memory multi-index db with cooperative runtime

История Tarantool



in-memory persistent cache

2008

in-memory persistent cache with replication

in-memory key/value database

2010

in-memory multi-index database (ACID)

in-memory multi-index db with lua functions

2012

in-memory multi-index db with cooperative runtime

hybrid multi-index db with sql and lua app server

2016

История Tarantool



in-memory persistent cache

2008

in-memory persistent cache with replication

in-memory key/value database

2010

in-memory multi-index database (ACID)

in-memory multi-index db with lua functions

2012

in-memory multi-index db with cooperative runtime

hybrid multi-index db with sql and lua app server *

2016



Tarantool

42

Яндекс



* In-memory computing platform. With SQL/NoSQL In-memory database + Appliances

История Tarantool



in-memory **persistent** cache

2008

in-memory persistent cache with **replication**

in-memory key/value database

2010

in-memory **multi-index** database (**ACID**)

in-memory multi-index db with lua functions

2012

in-memory multi-index db with cooperative runtime

hybrid multi-index db with sql and lua app server *

2016



Tarantool

43

Яндекс



* In-memory computing platform. With SQL/NoSQL In-memory database + Appliances

ACID

Atomicity

Cooperative, 1 thread

Consistency

Coop monopoly of txn's

Isolation

No parallel txn's

Durability

Write Ahead Log

Cooperative vs Multithread

RAM: 0.5 М - 2.5 М чтений/s блоками по 4kb

1 чтение ~ 1 мкс => 1M RPS

Запись не требует блокировок

Основной код базы — не требует потокобезопасности

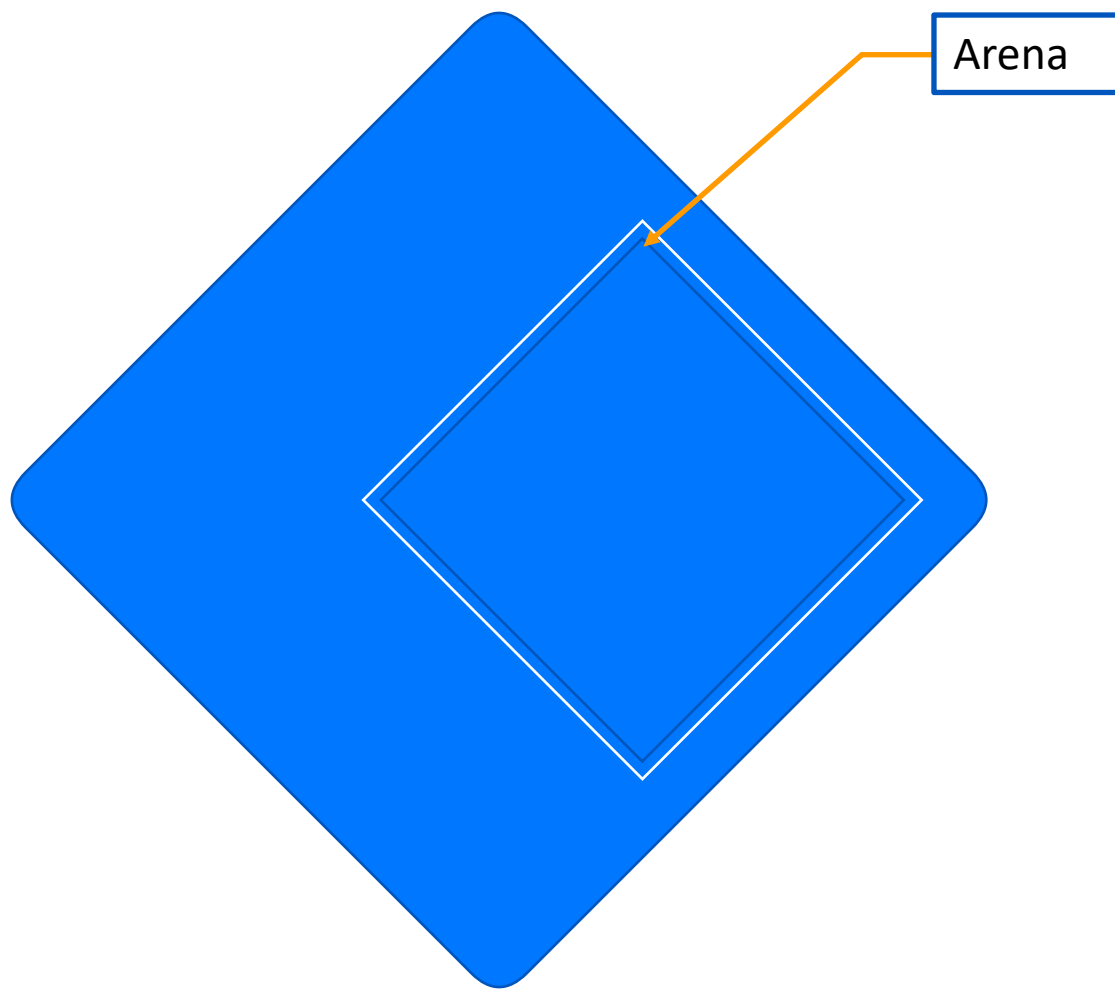
Как устроен Tarantool?

- 100% данных находятся в памяти
- Доступ к данным из одного потока
- Доступ к данным только через индекс
- Изменения пишутся во Write Ahead Log (WAL)
- WAL реплицируется
- Периодически сохраняются консистентный Snapshot

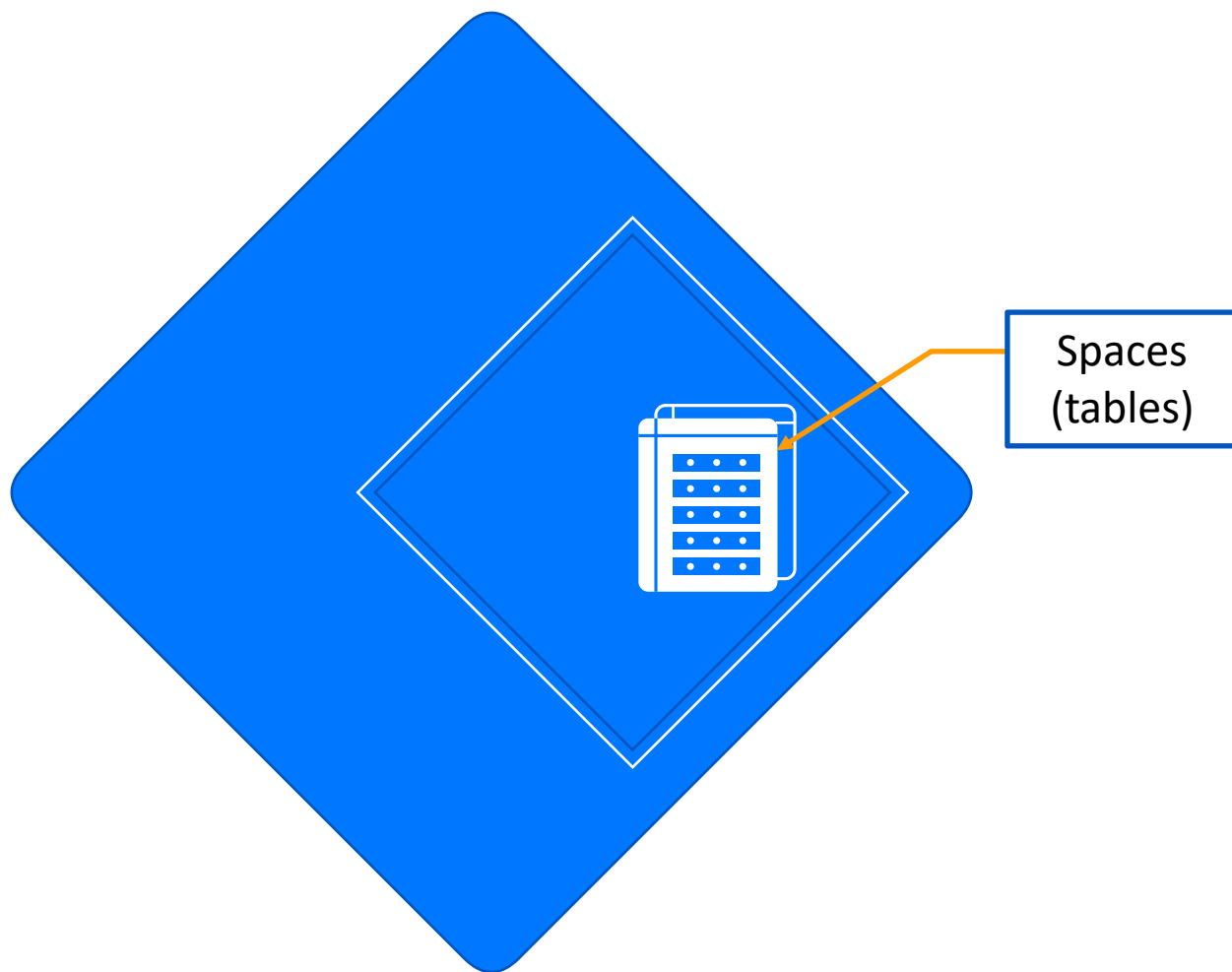


TX
(transaction thread)

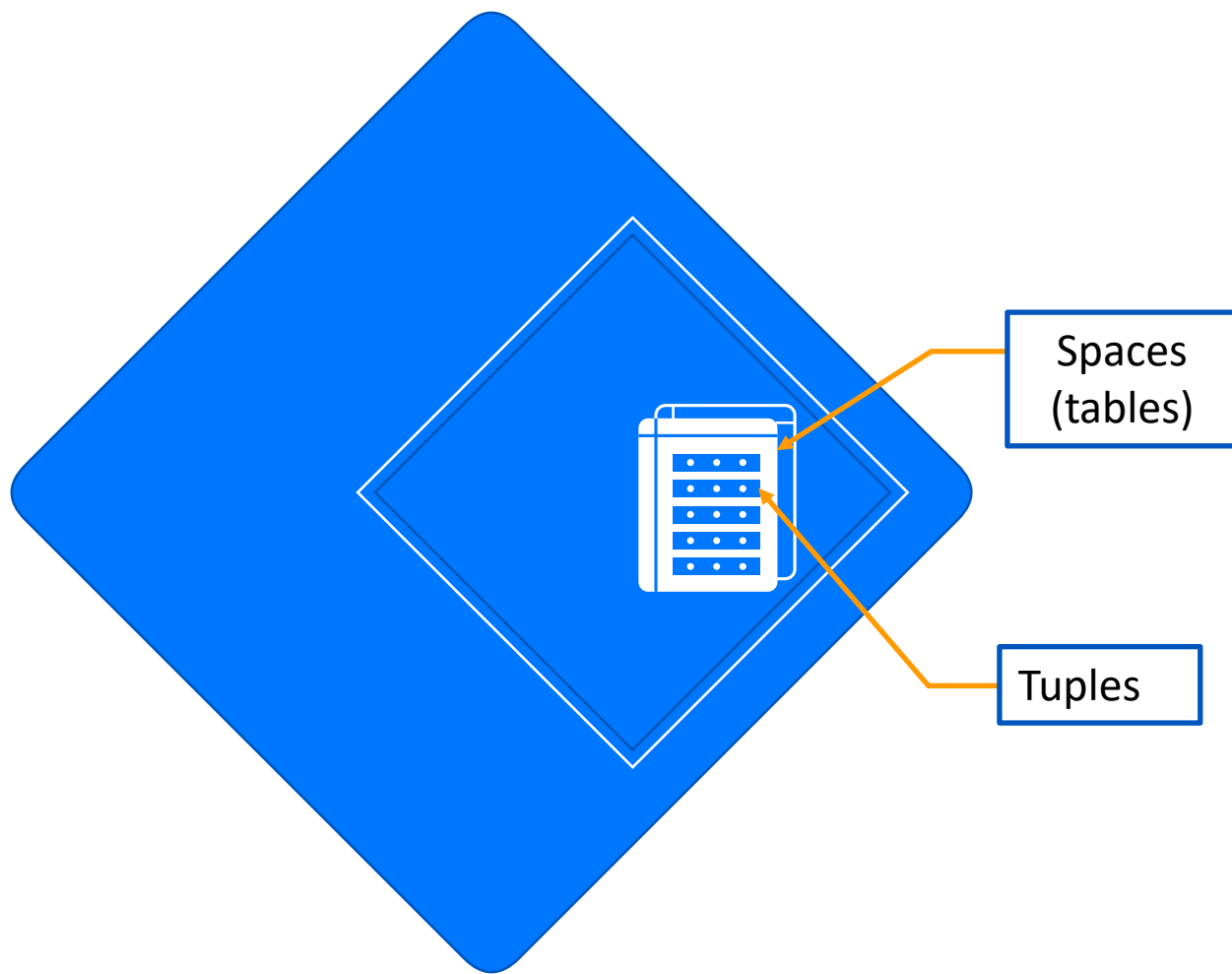
- Монопольный доступ



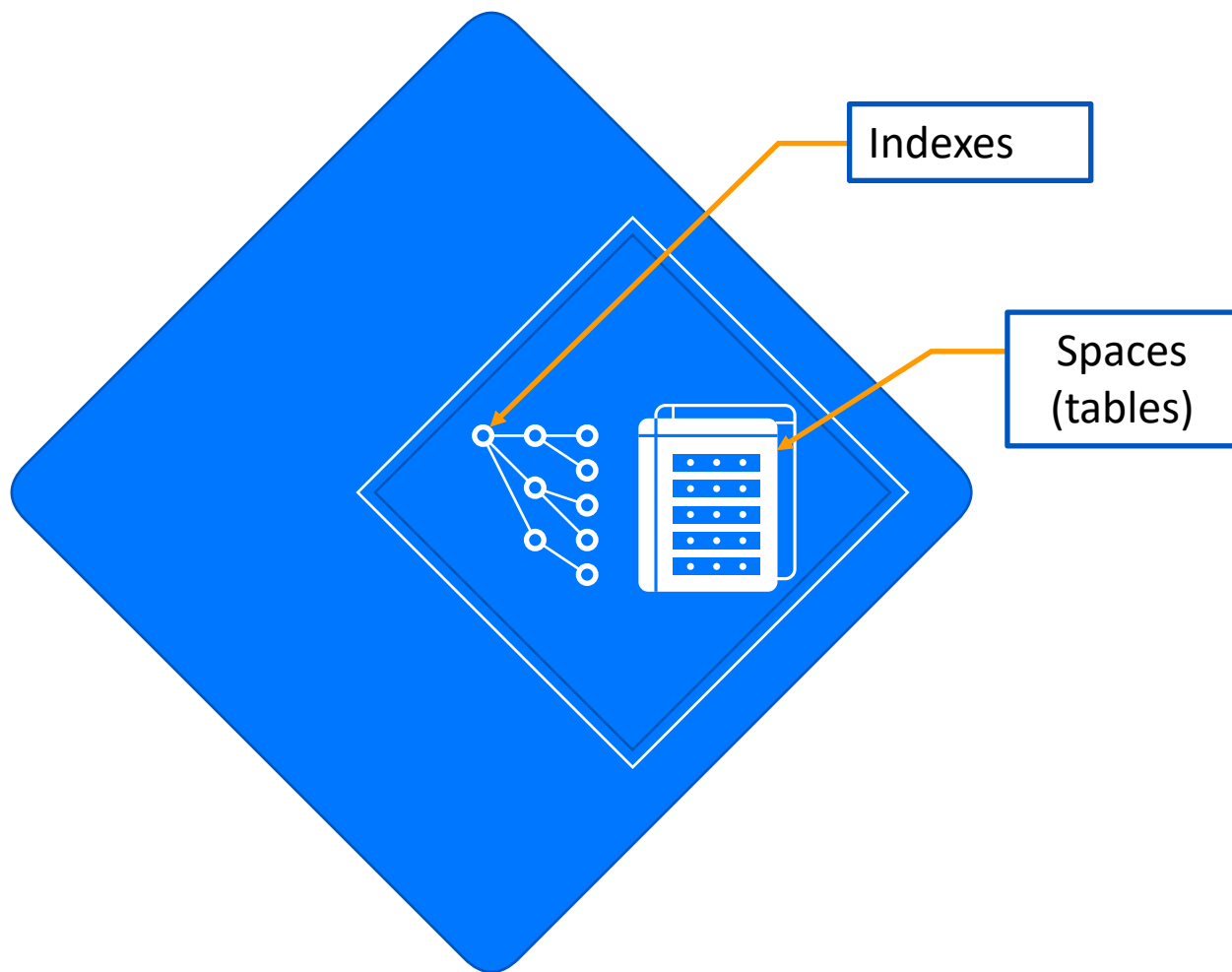
- Монопольный доступ
- Данные в памяти



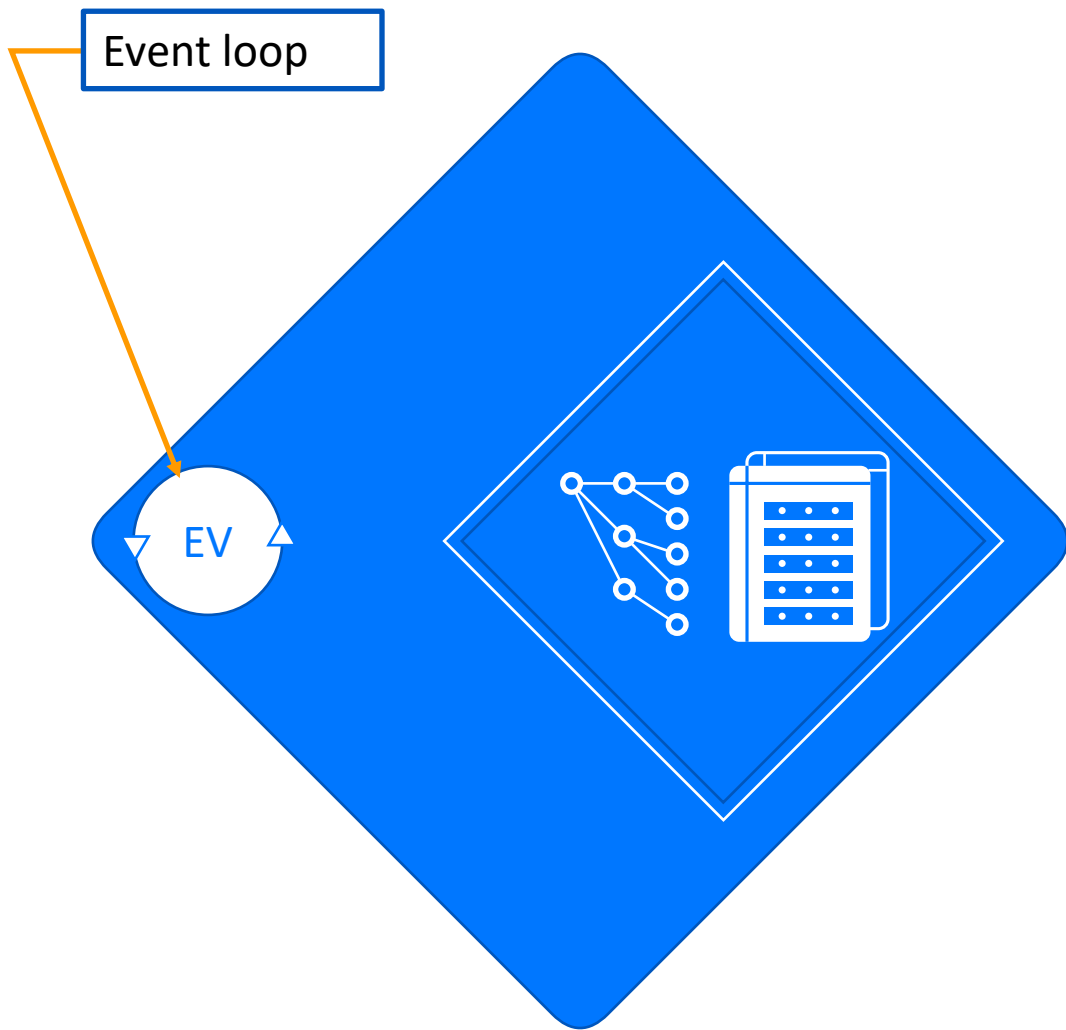
- Монопольный доступ
- Данные в памяти



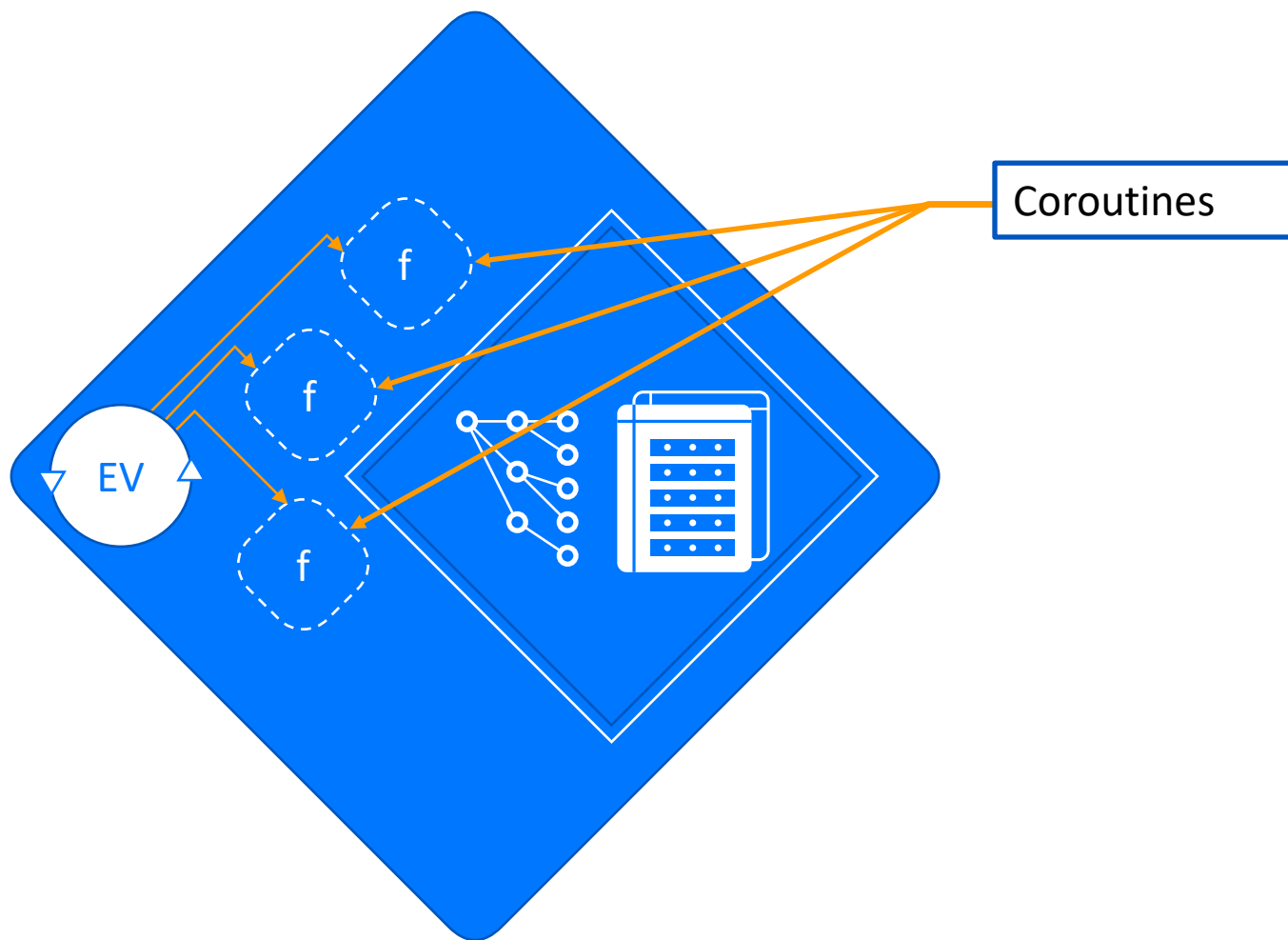
- Монопольный доступ
- Данные в памяти



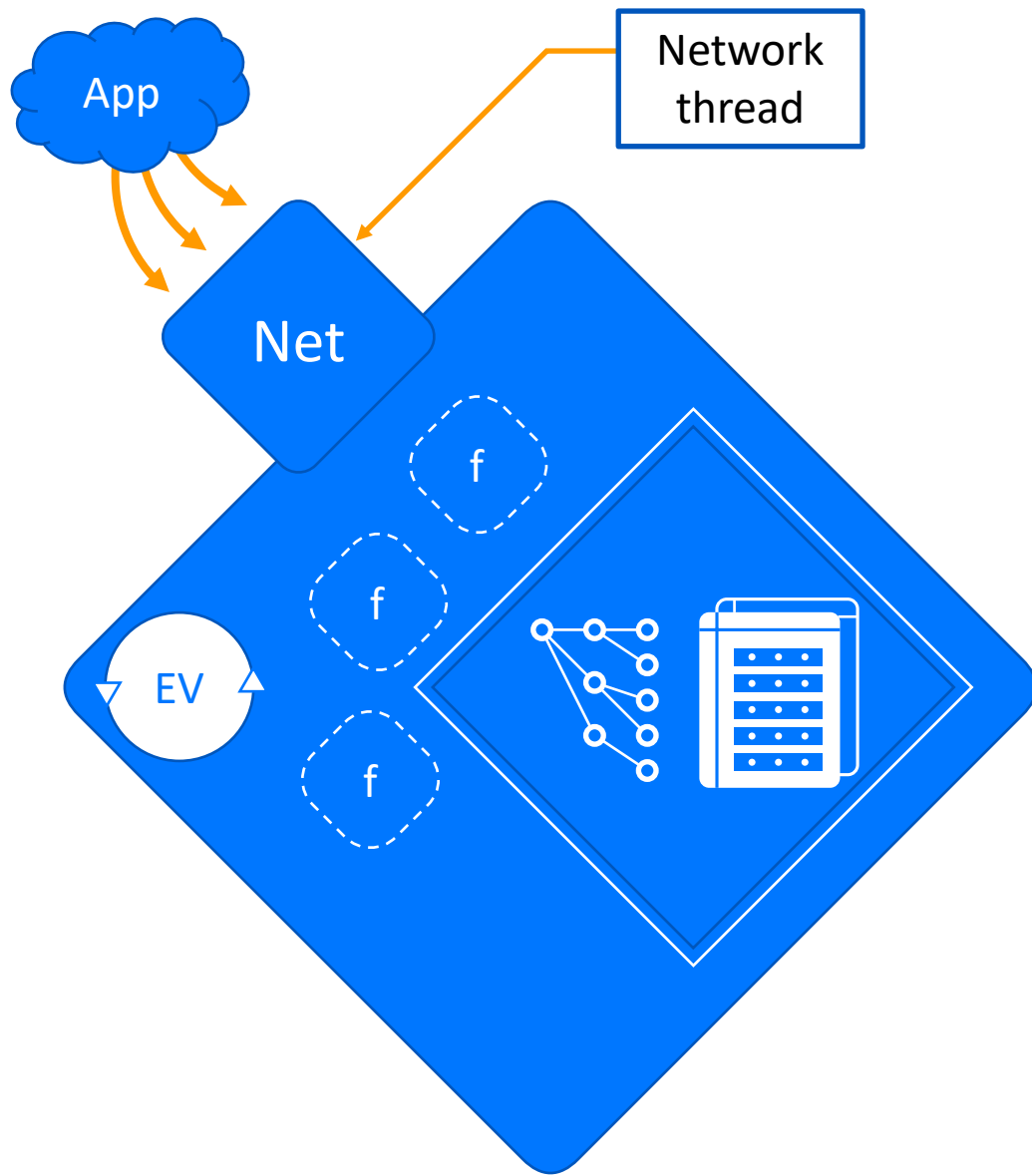
- Монопольный доступ
- Данные в памяти



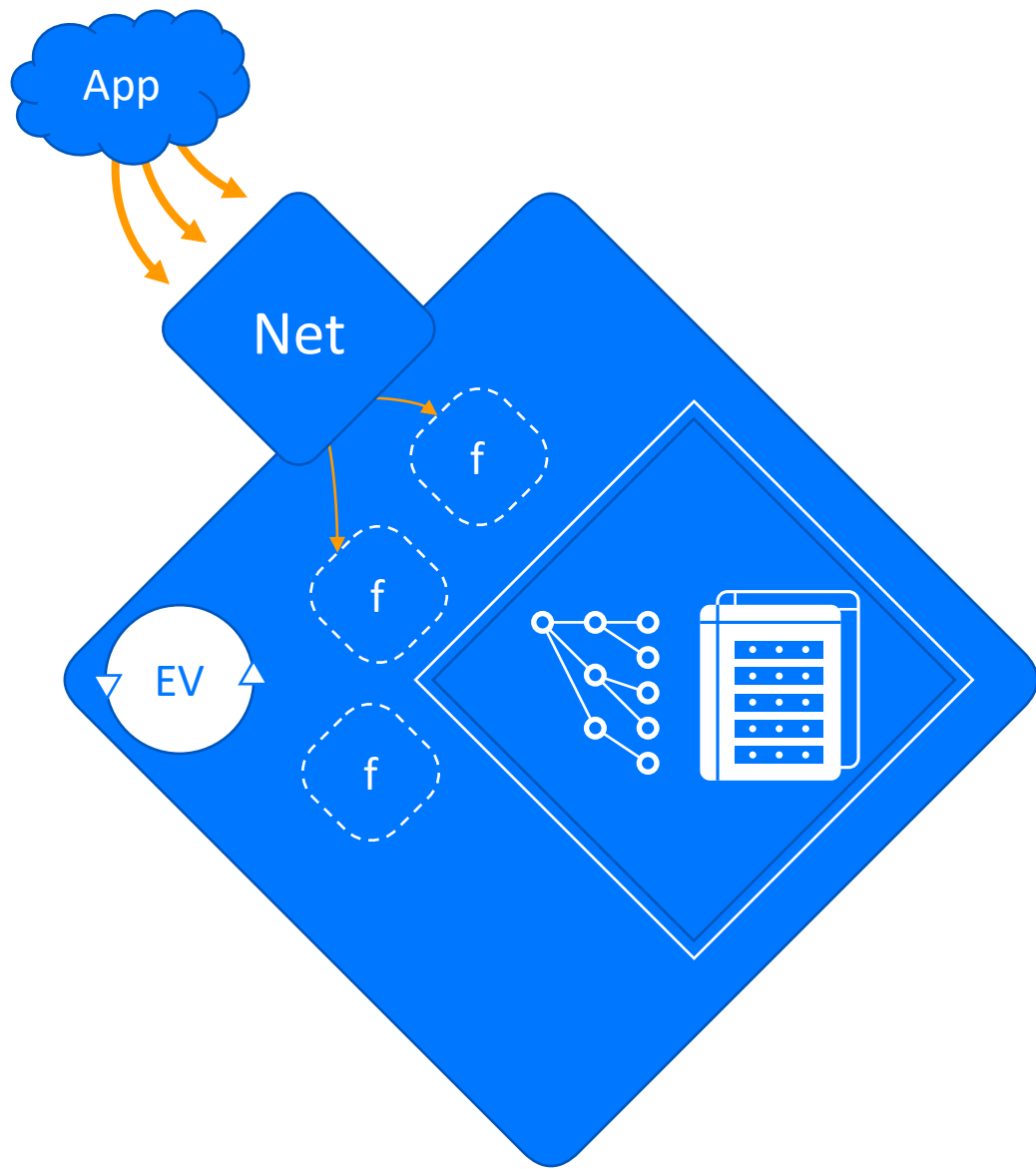
- Данные в памяти
- Монопольный доступ
- Событийный цикл



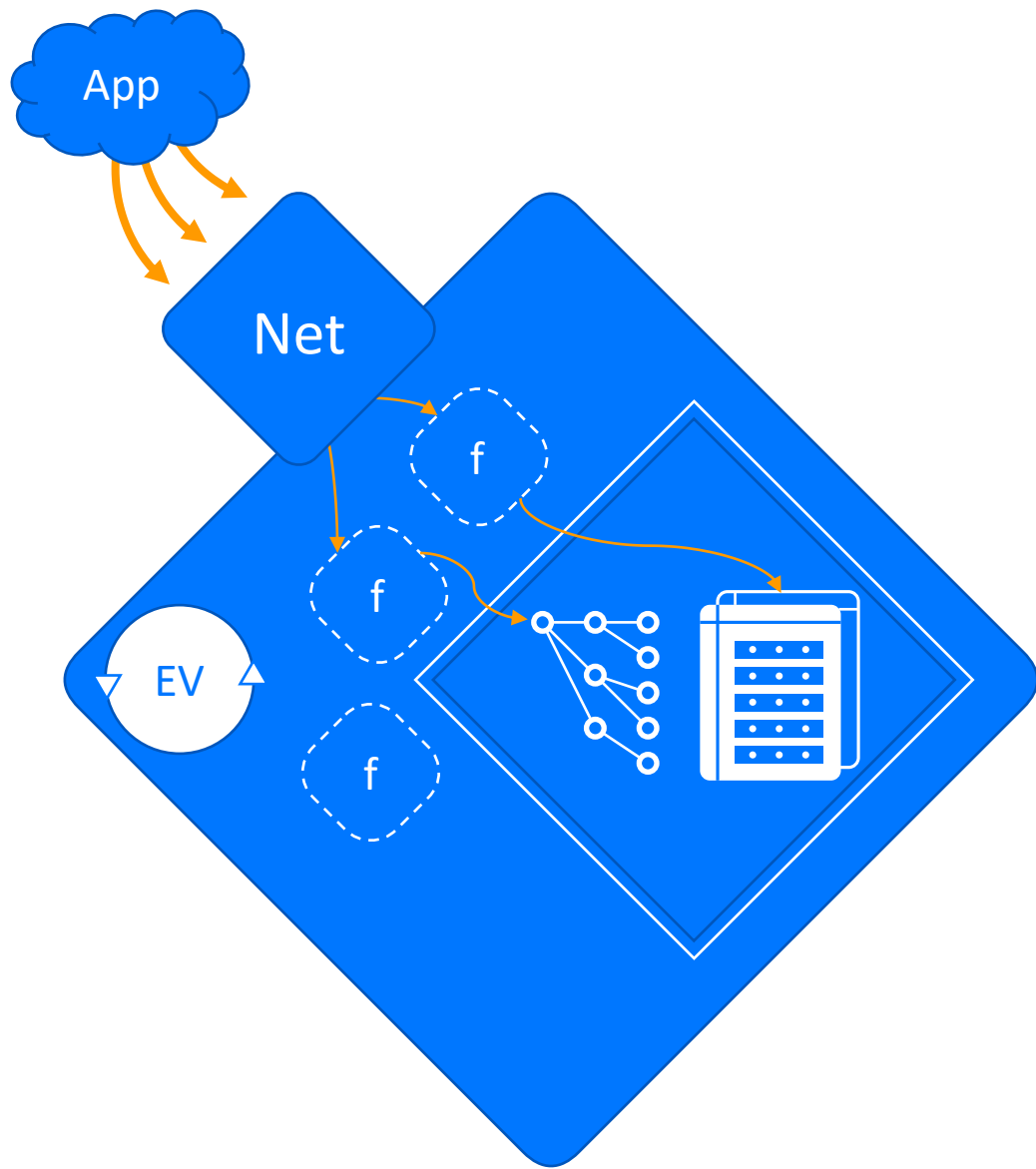
- Данные в памяти
- Монопольный доступ
- Событийный цикл
- Кооперативная многозадачность



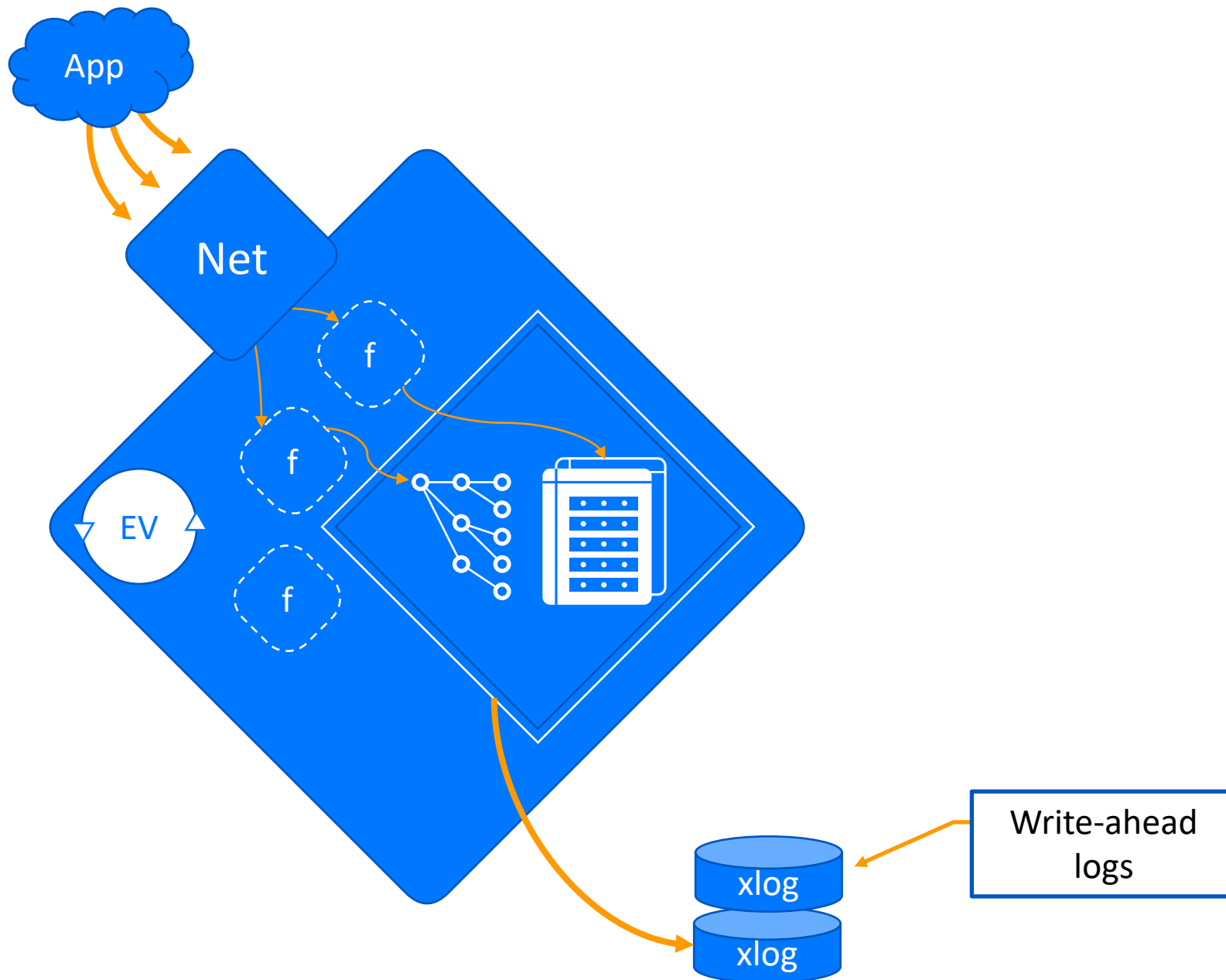
- Данные в памяти
- Монопольный доступ
- Событийный цикл
- Кооперативная многозадачность
- Отдельные потоки для сети



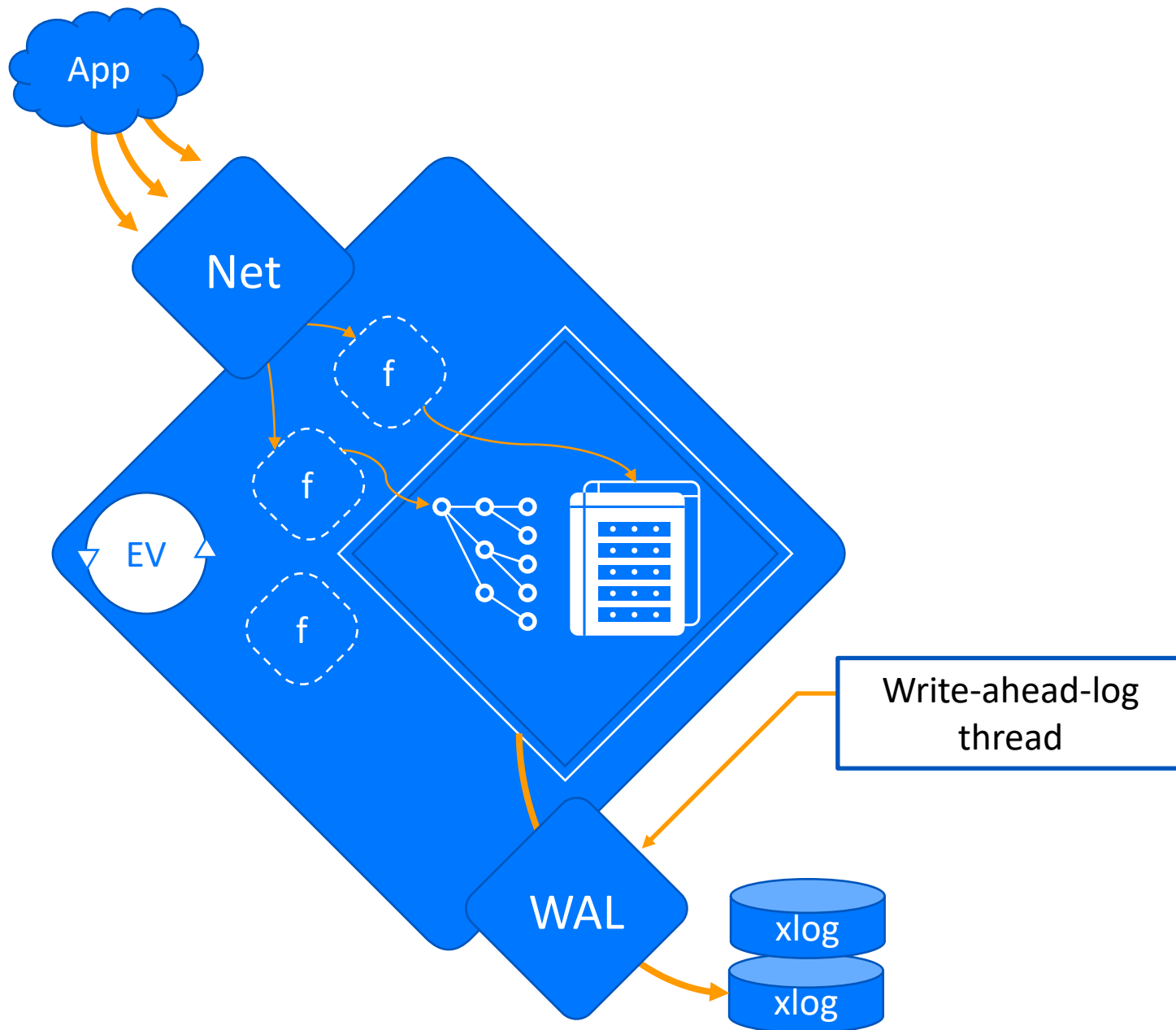
- Данные в памяти
- Монопольный доступ
- Событийный цикл
- Кооперативная многозадачность
- Отдельные потоки для сети



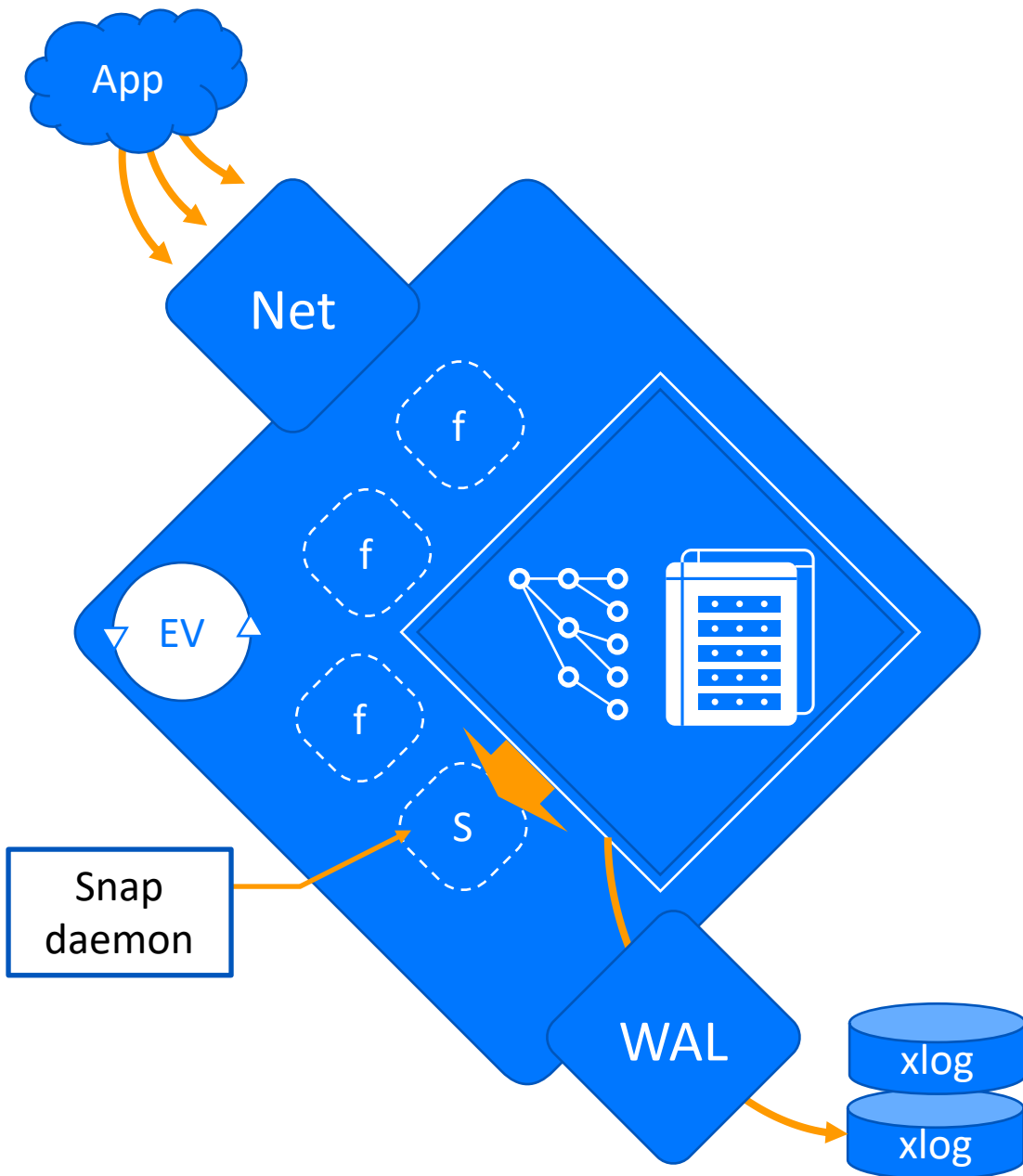
- Данные в памяти
- Монопольный доступ
- Событийный цикл
- Кооперативная многозадачность
- Отдельные потоки для сети



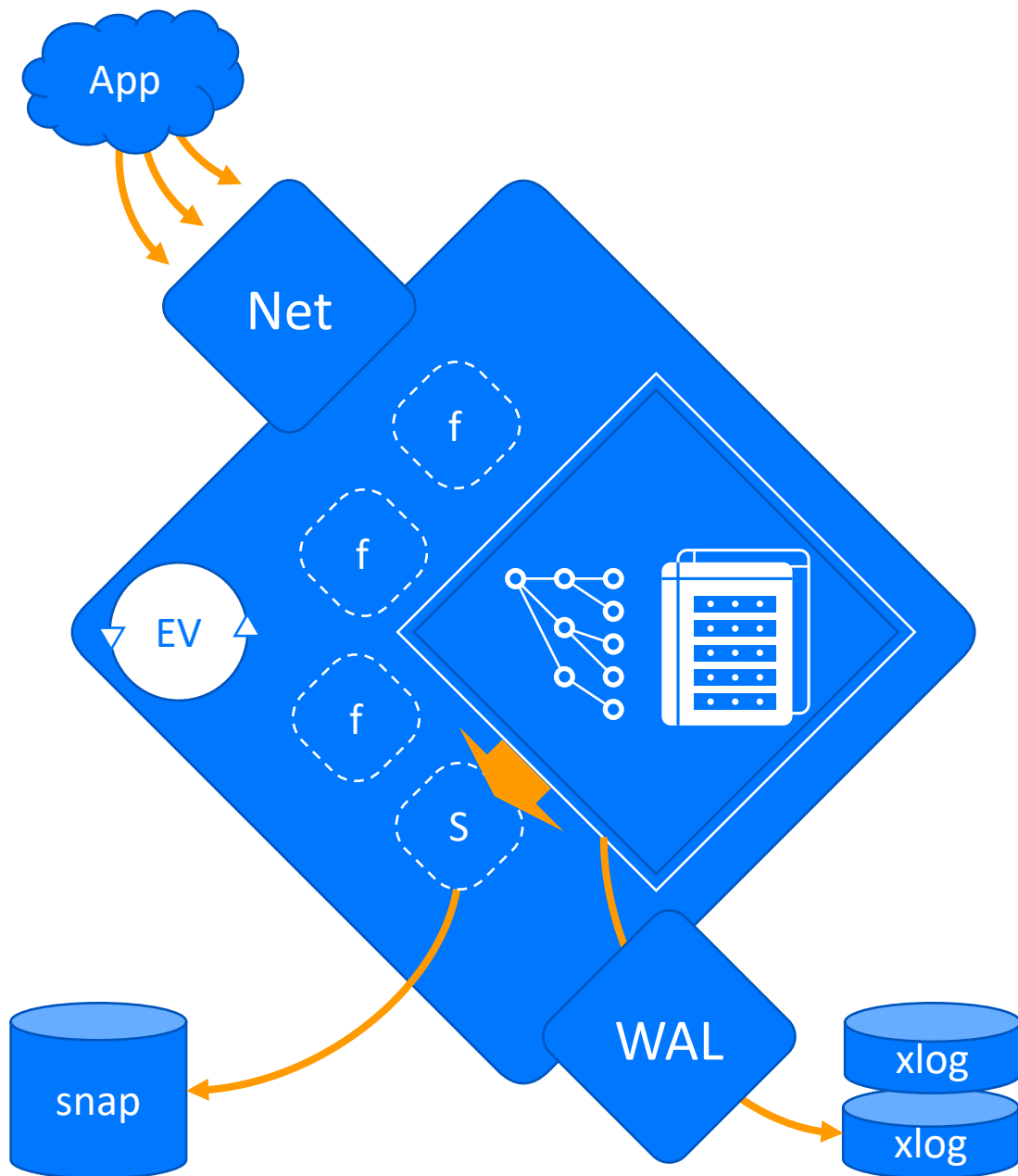
- Данные в памяти
- Монопольный доступ
- Событийный цикл
- Кооперативная многозадачность
- Отдельные потоки для сети
- Изменения во Write-Ahead-Log



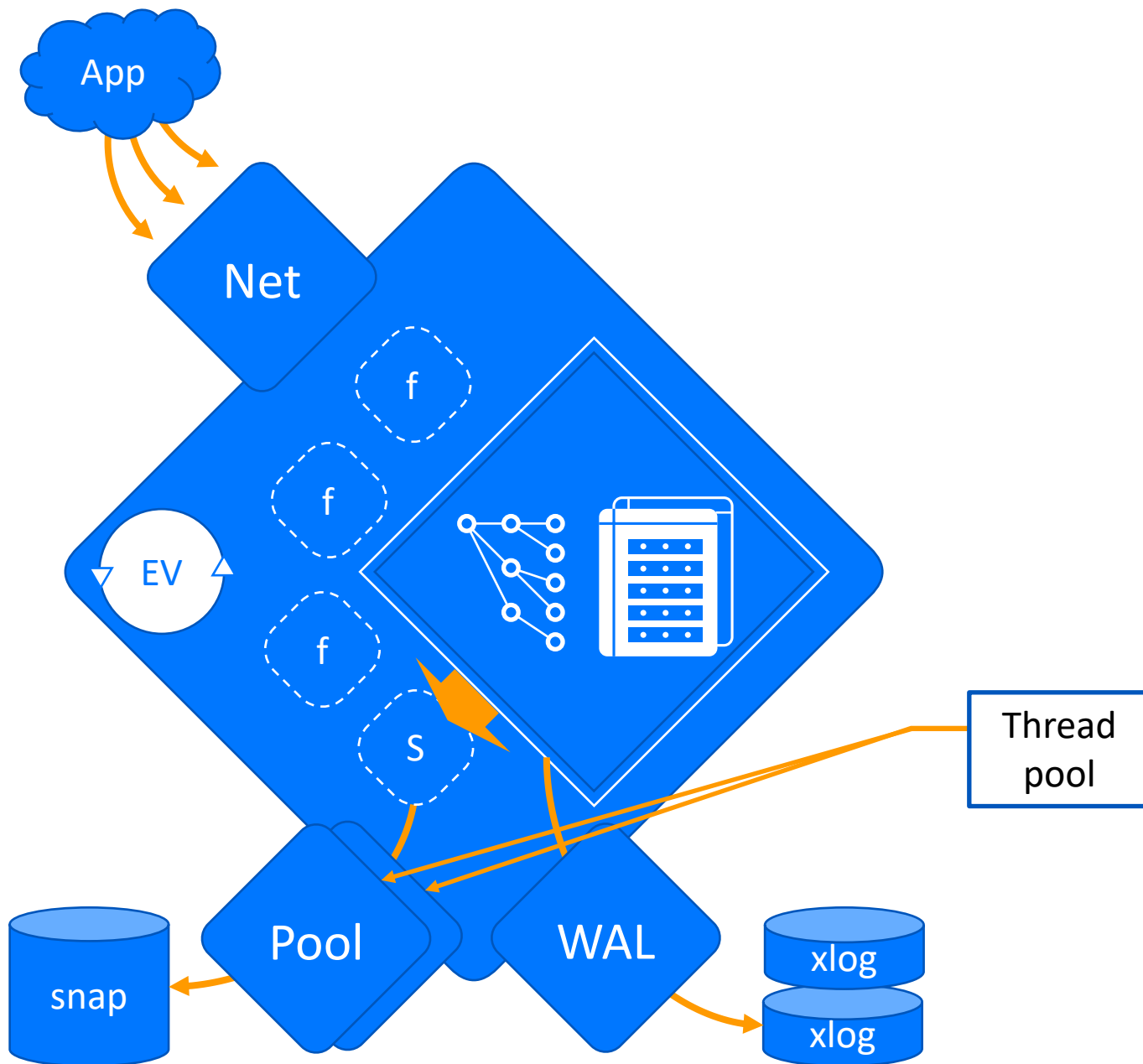
- Данные в памяти
- Монопольный доступ
- Событийный цикл
- Кооперативная многозадачность
- Отдельные потоки для сети
- Изменения во Write-Ahead-Log



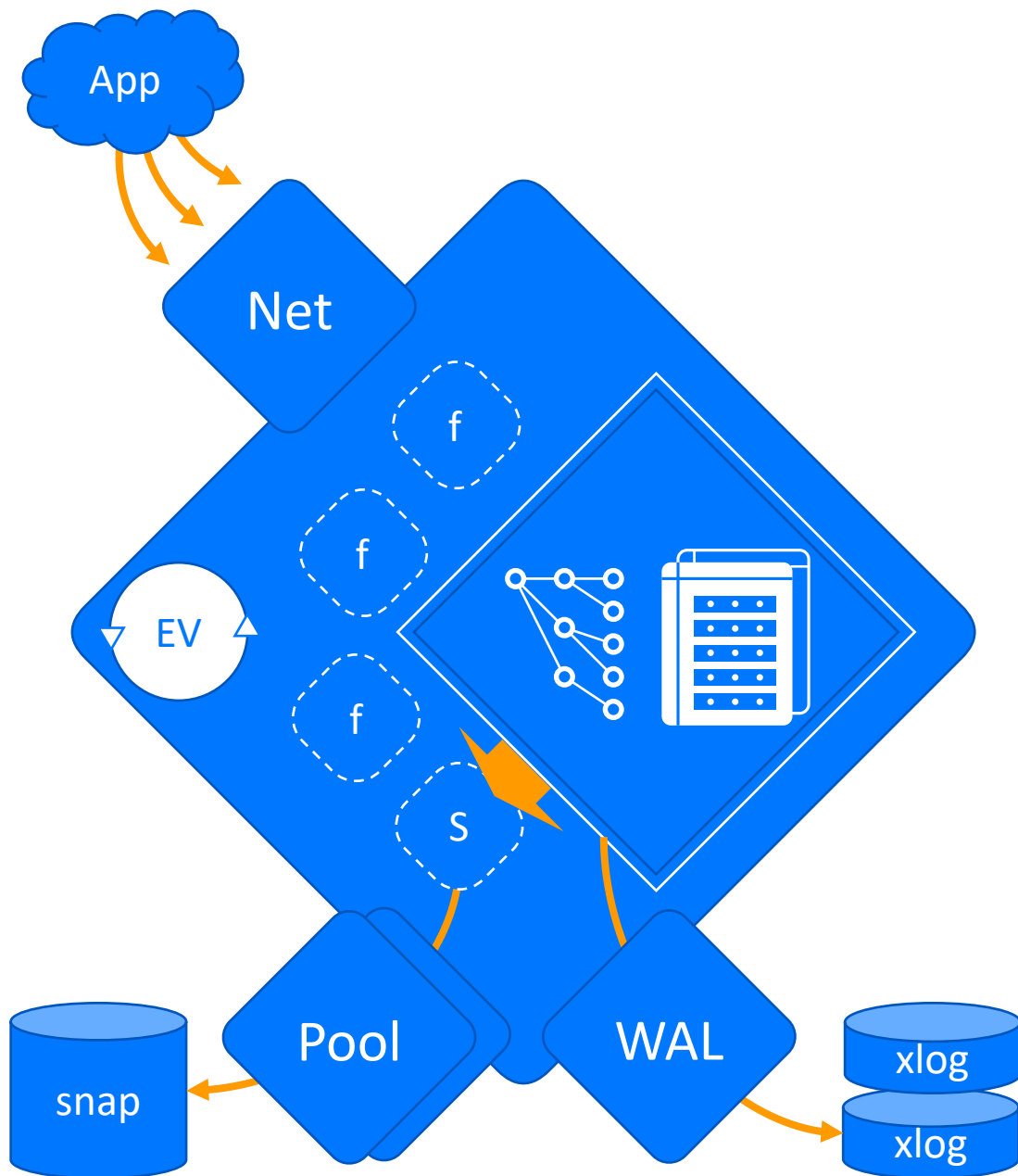
- Данные в памяти
- Монопольный доступ
- Событийный цикл
- Кооперативная многозадачность
- Отдельные потоки для сети
- Изменения во Write-Ahead-Log
- Консистентный снимок



- Данные в памяти
- Монопольный доступ
- Событийный цикл
- Кооперативная многозадачность
- Отдельные потоки для сети
- Изменения во Write-Ahead-Log
- Консистентный снапшот



- Данные в памяти
- Монопольный доступ
- Событийный цикл
- Кооперативная многозадачность
- Отдельные потоки для сети
- Изменения во Write-Ahead-Log
- Консистентный снимок
- Отдельные потоки для диска



- Данные в памяти
- Монопольный доступ
- Событийный цикл
- Кооперативная многозадачность
- Отдельные потоки для сети
- Изменения во Write-Ahead-Log
- Консистентный снапшот
- Отдельные потоки для диска

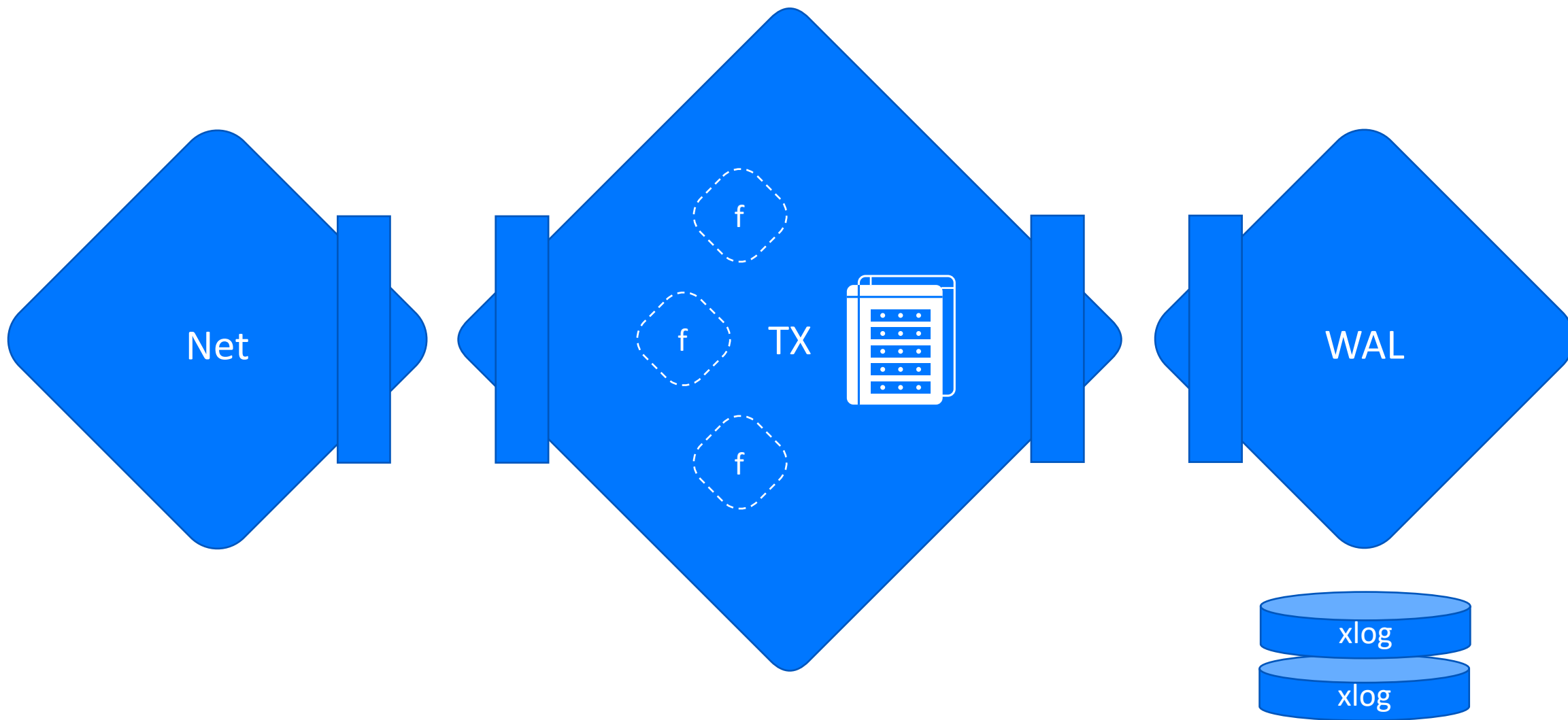
Откуда производительность?

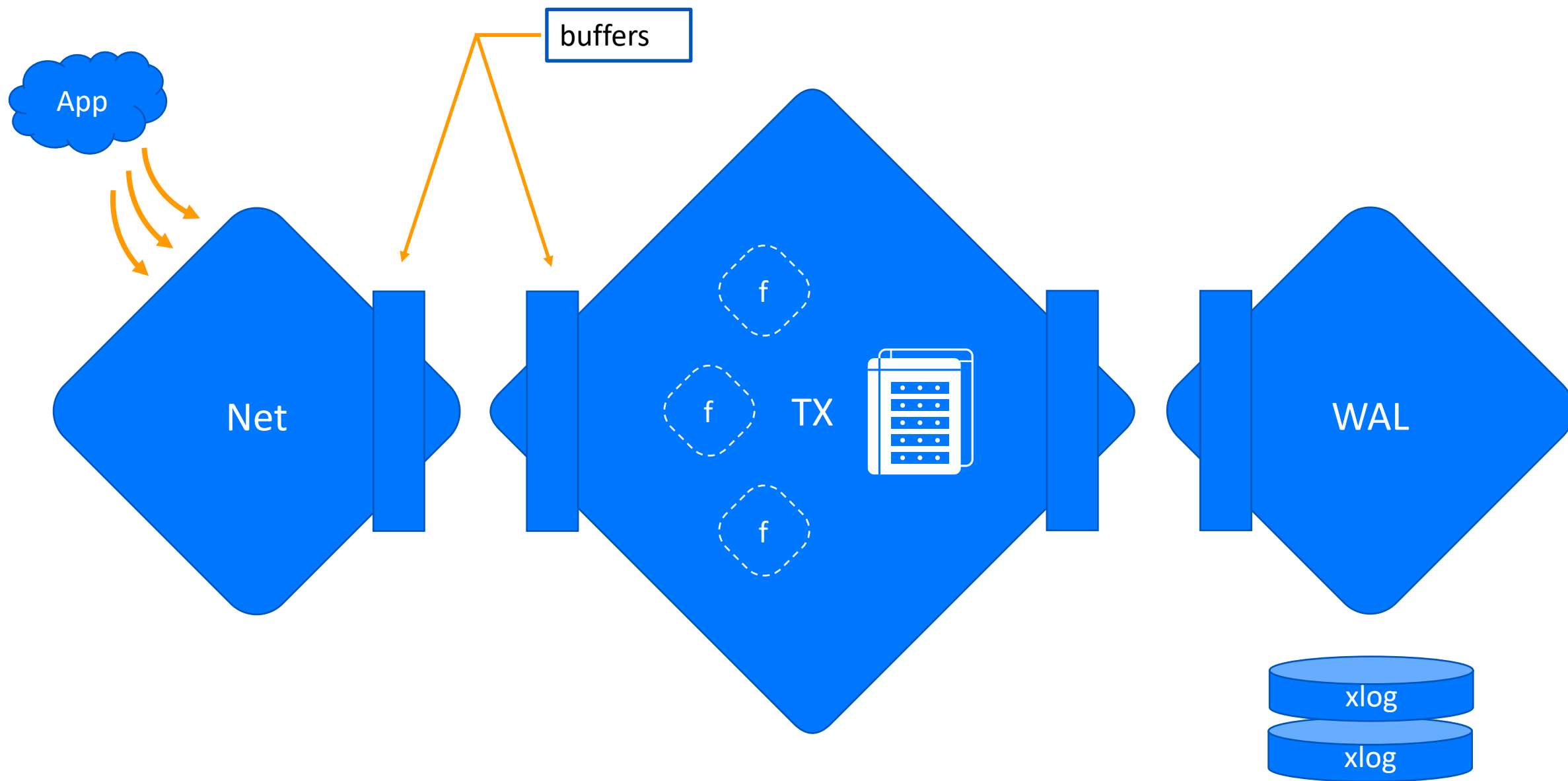
- Диск медленный
- Системные вызовы дорогие

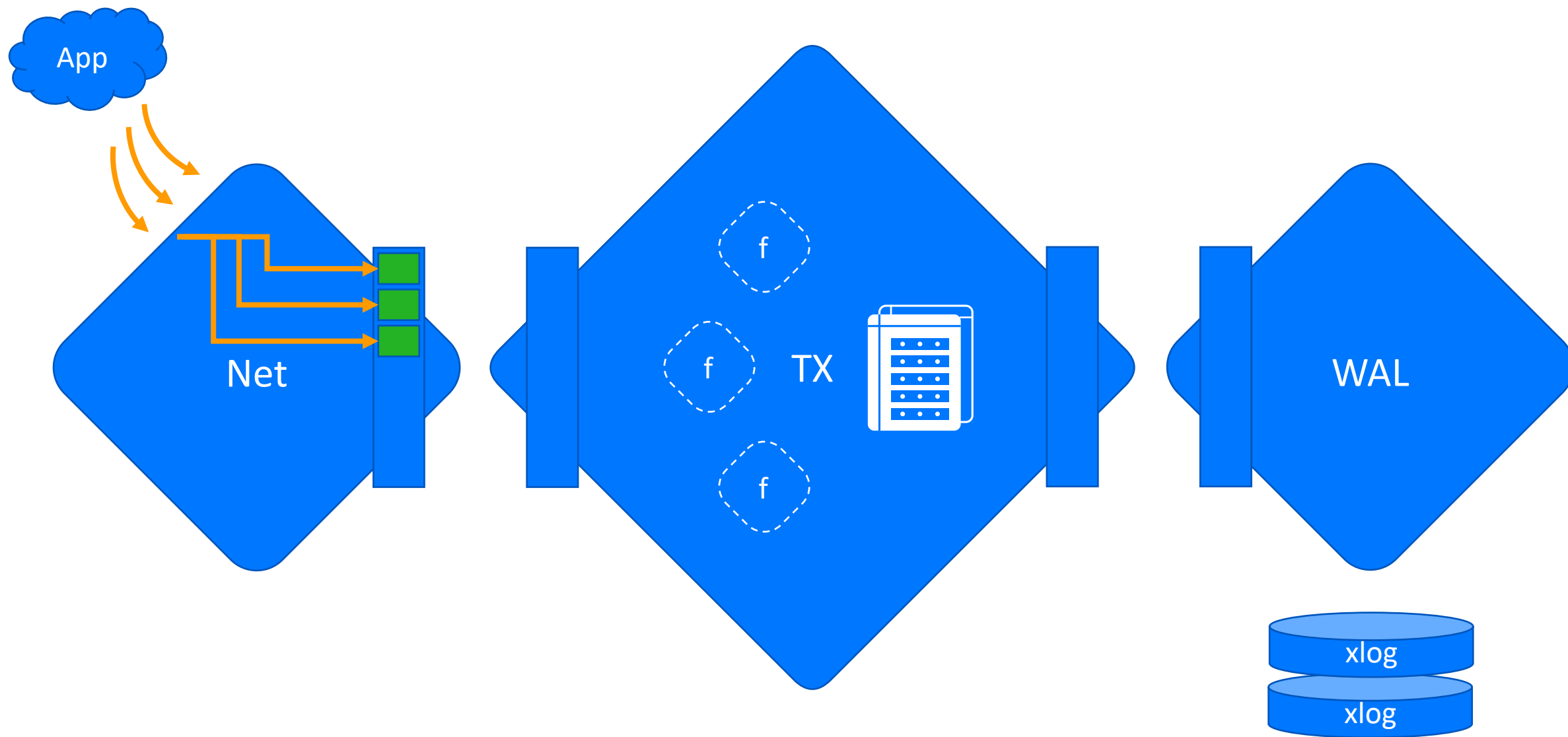
Откуда производительность?

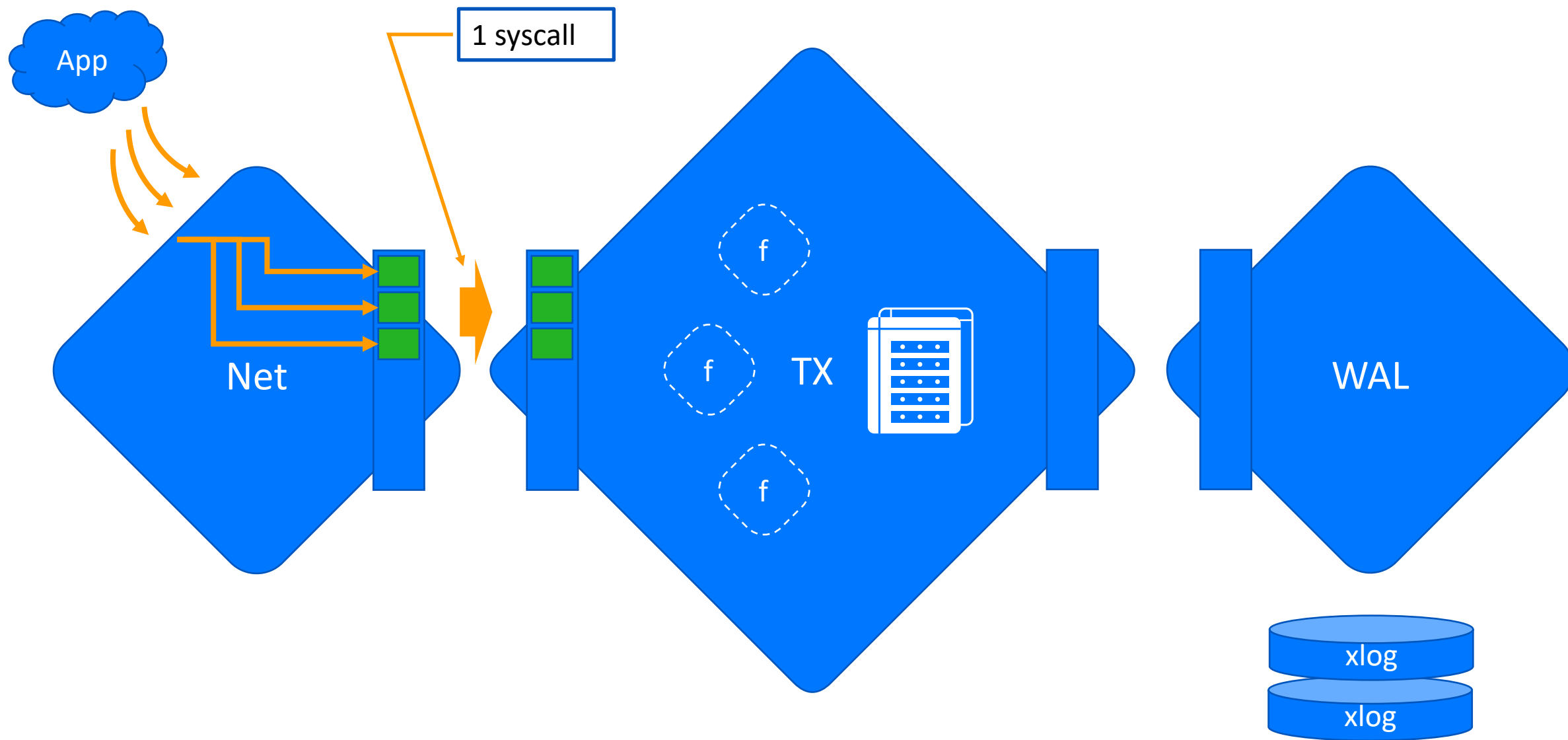
- Диск медленный
- Системные вызовы дорогие

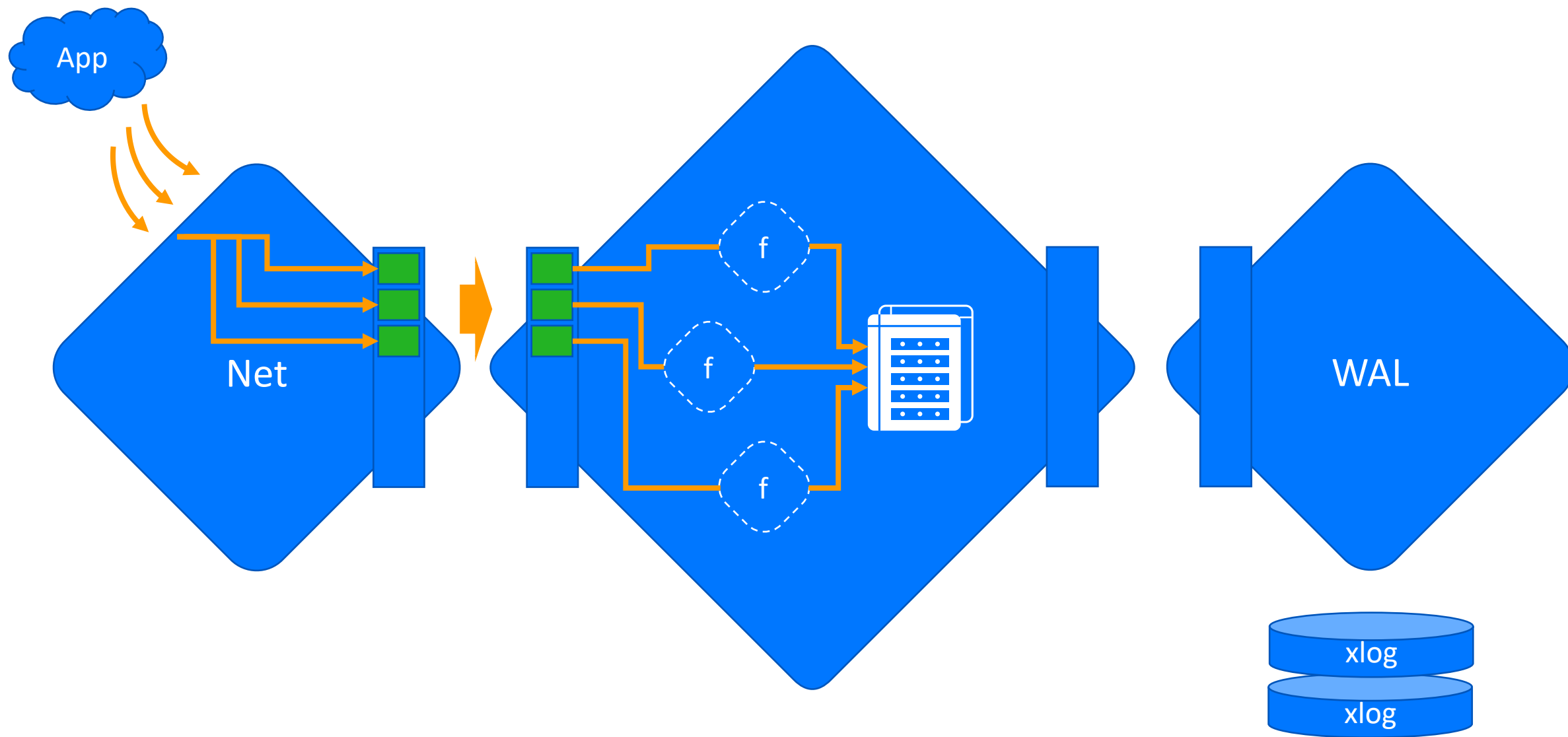
Батчинг!

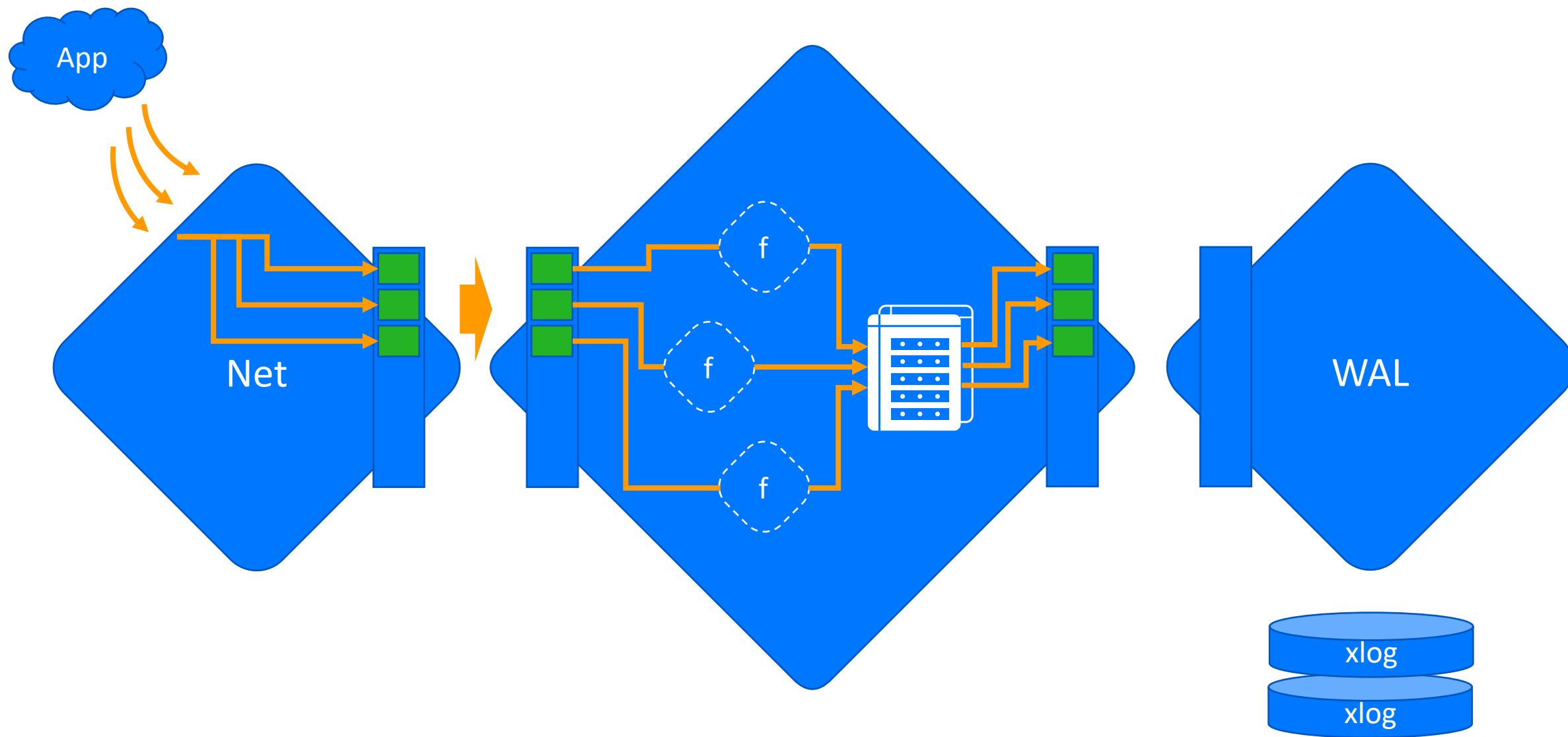


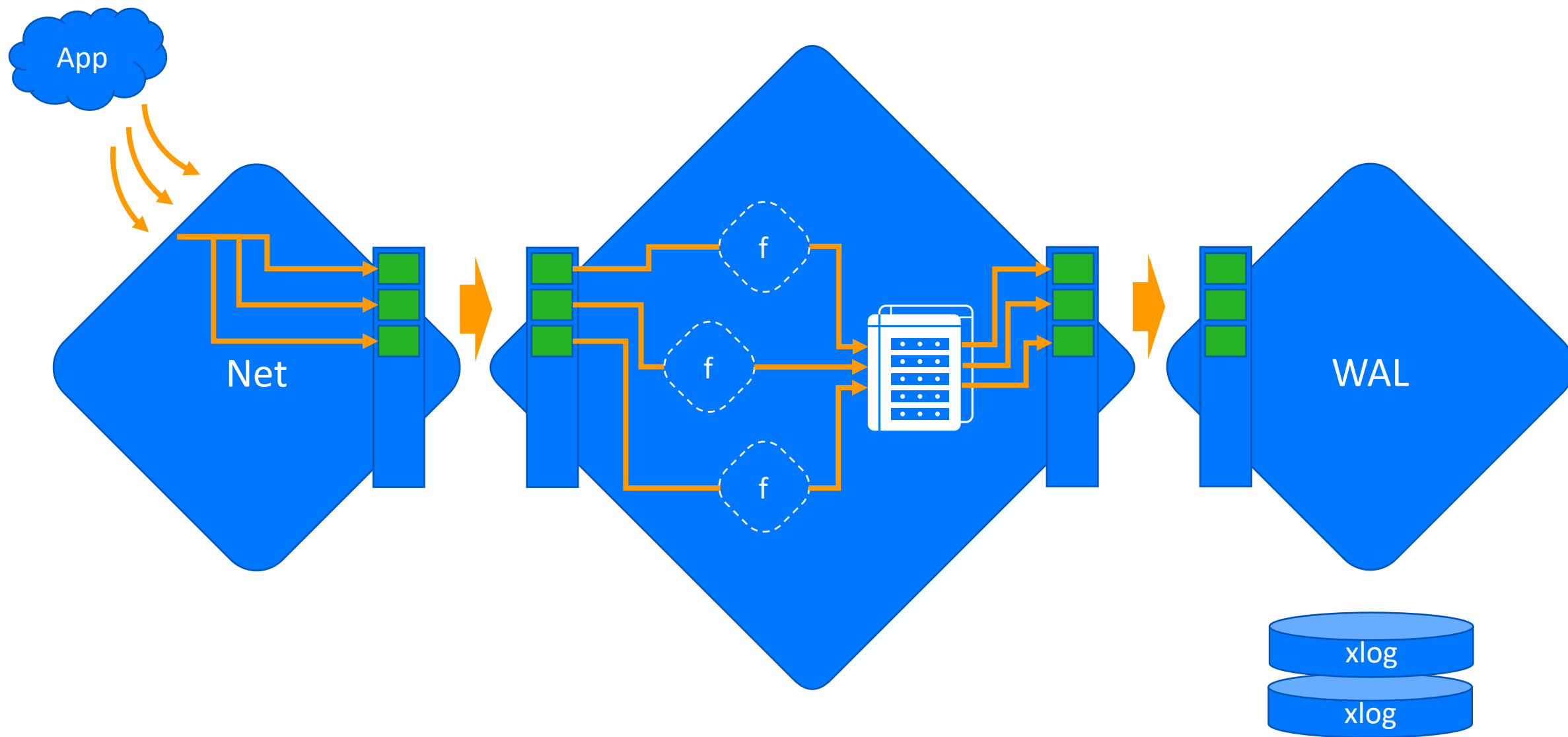


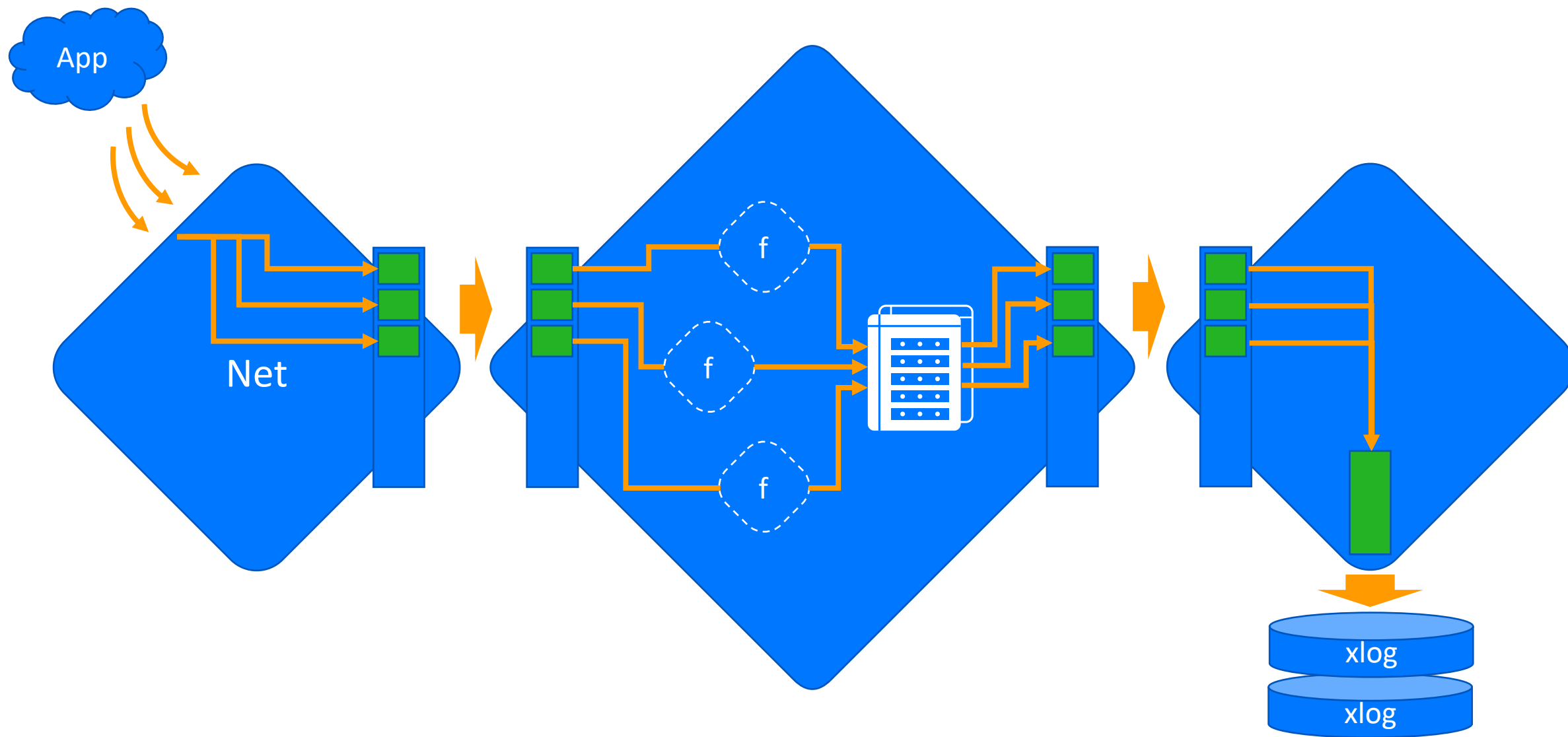


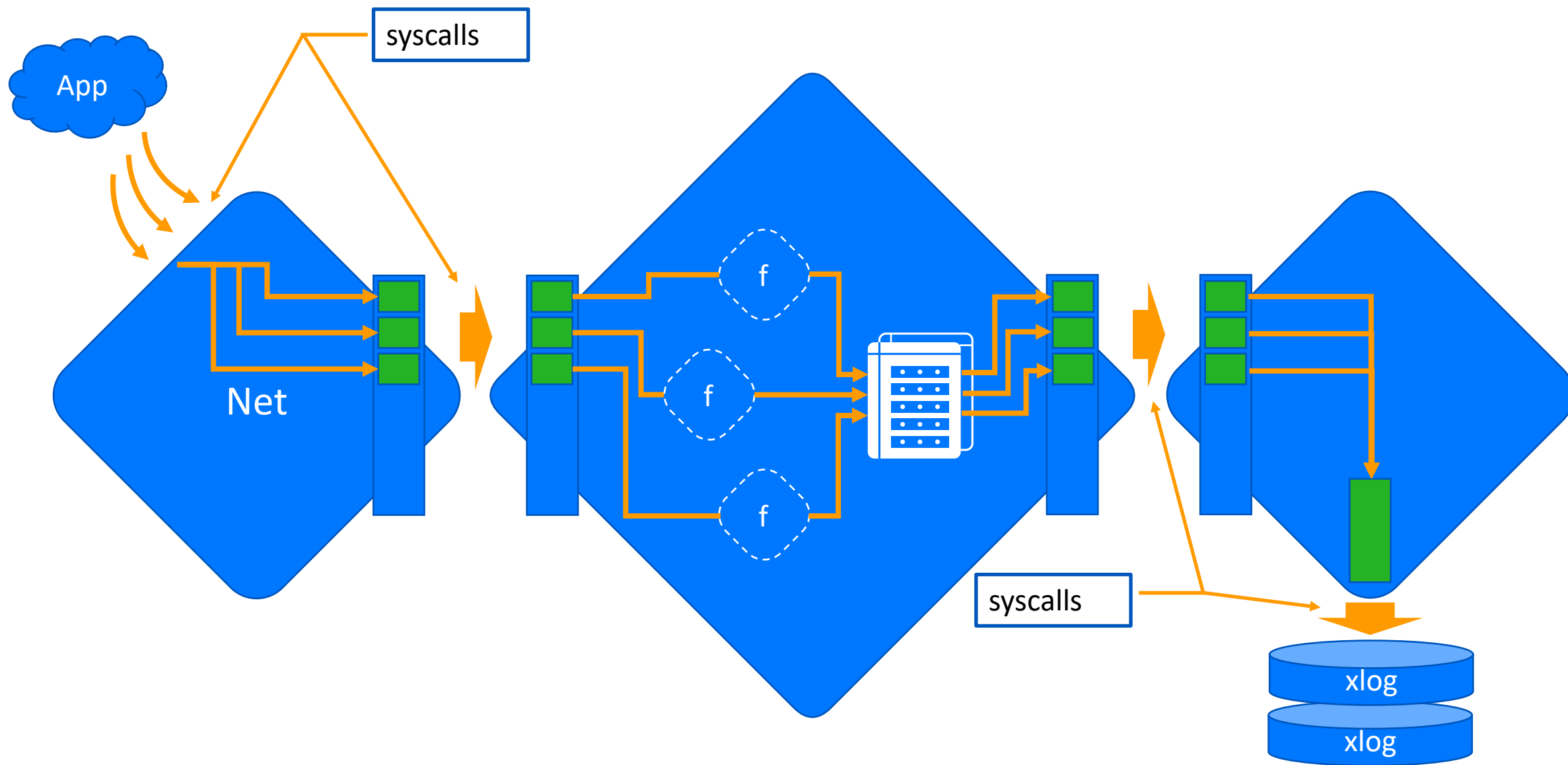


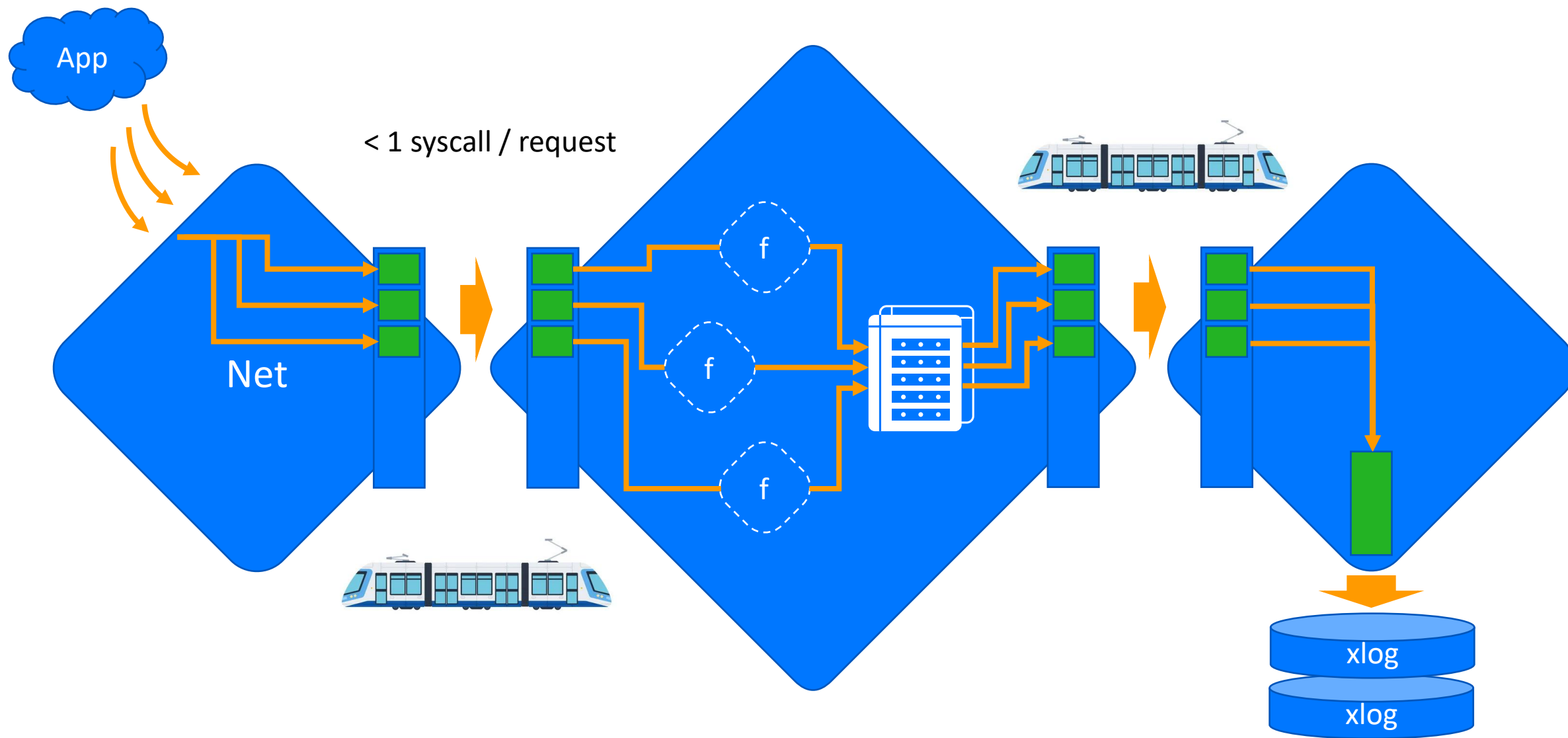


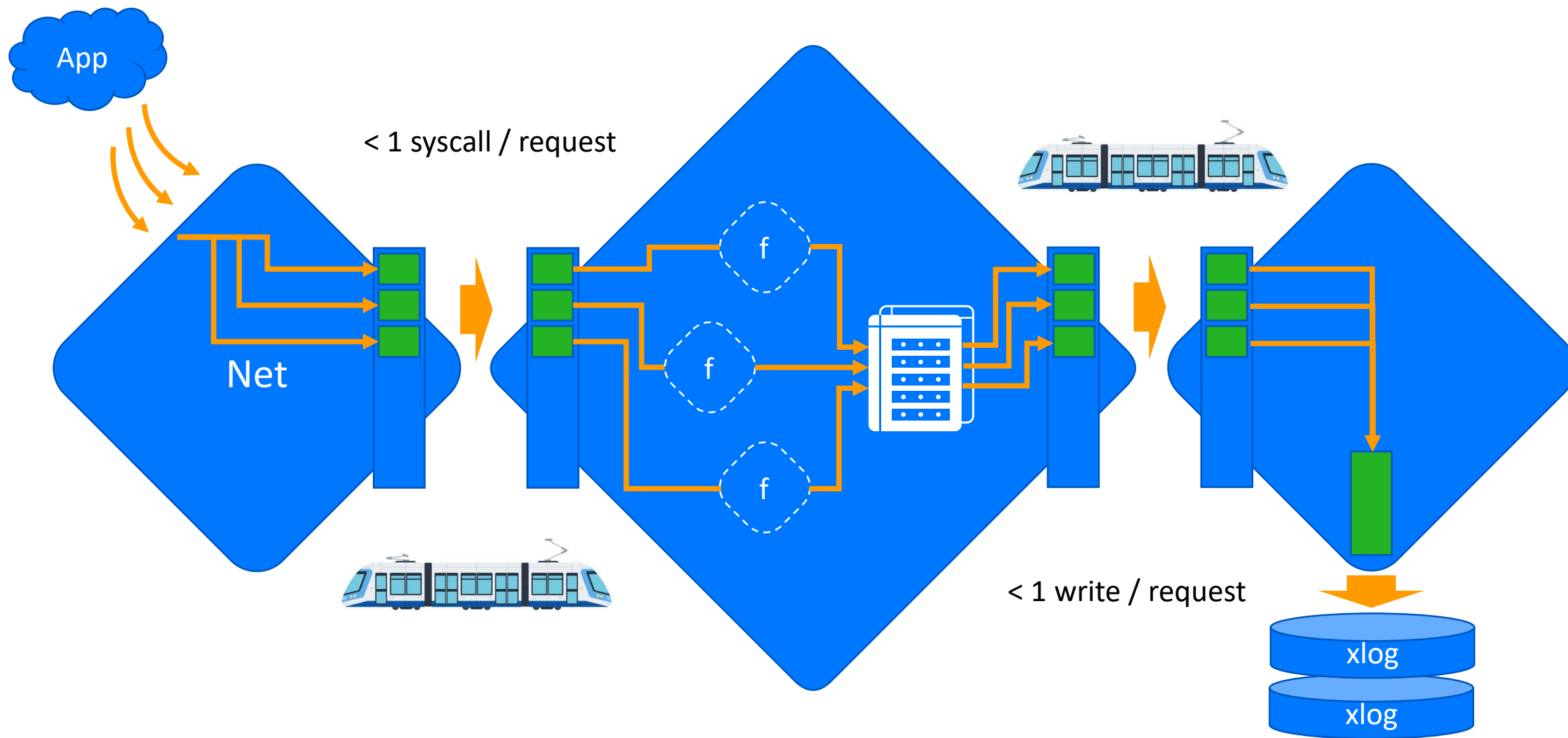


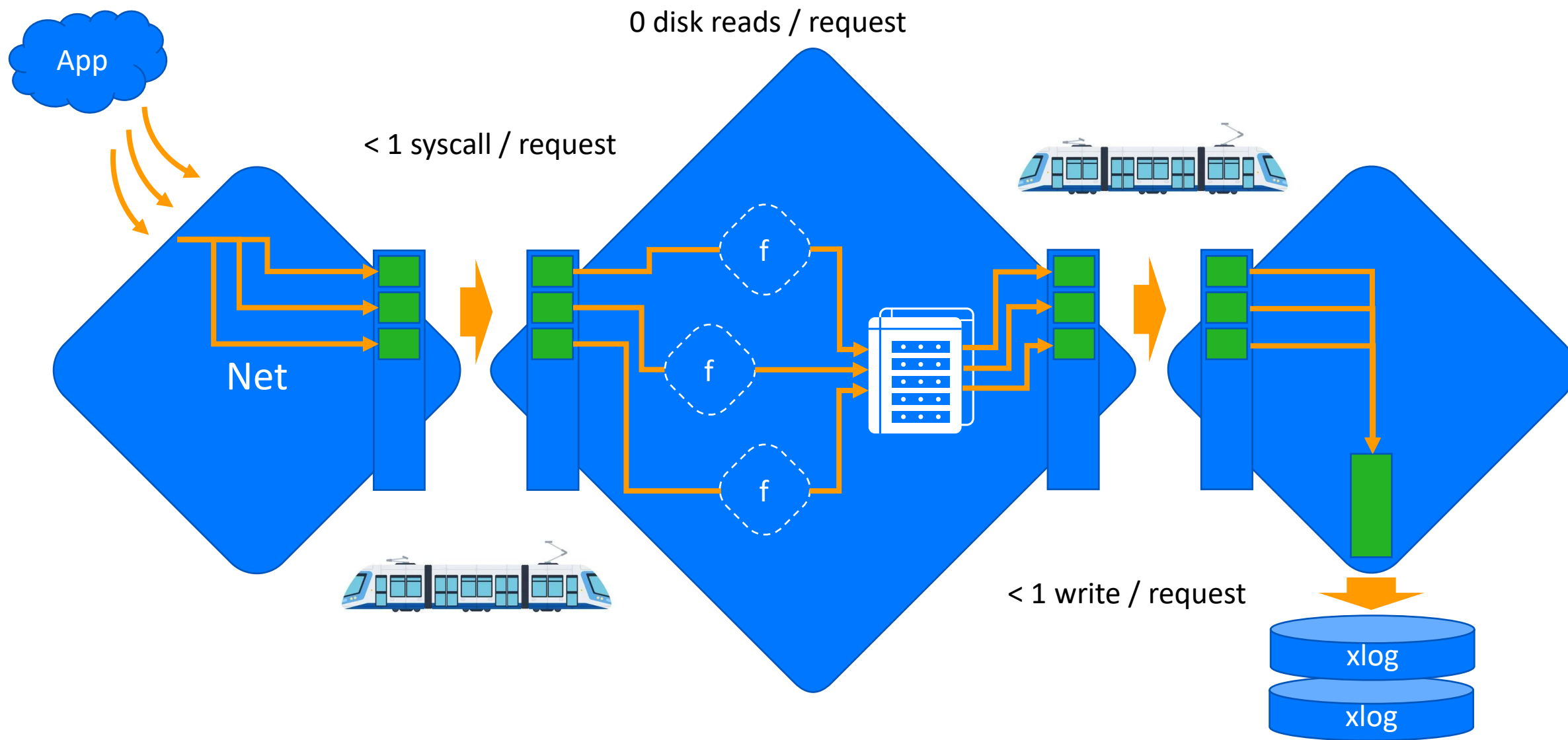












Сохранность данных в памяти

Запись «ждёт»

Пишем в память

Пишем на диск

Ждём подтверждения

Сохранность данных в памяти

Чтение «не ждёт»

Читаем из памяти

Гарантированно быстро

Запись «ждёт»

Пишем в память

Пишем на диск

Ждём подтверждения

Восстановление после перезапуска

- Чтение снимка
- Чтение xlog
- Построение индексов

Восстановление после перезапуска

- Чтение снимшота
- Чтение xlog
- Построение индексов

Долгий старт

03

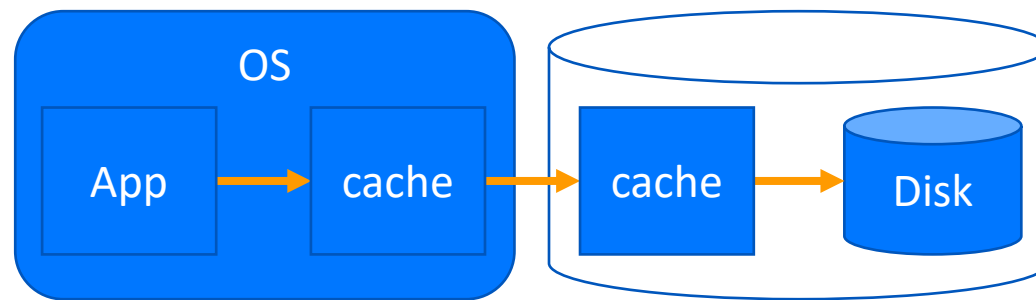
Обеспечение сохранности данных

Персистентность, диск и миф про fsync

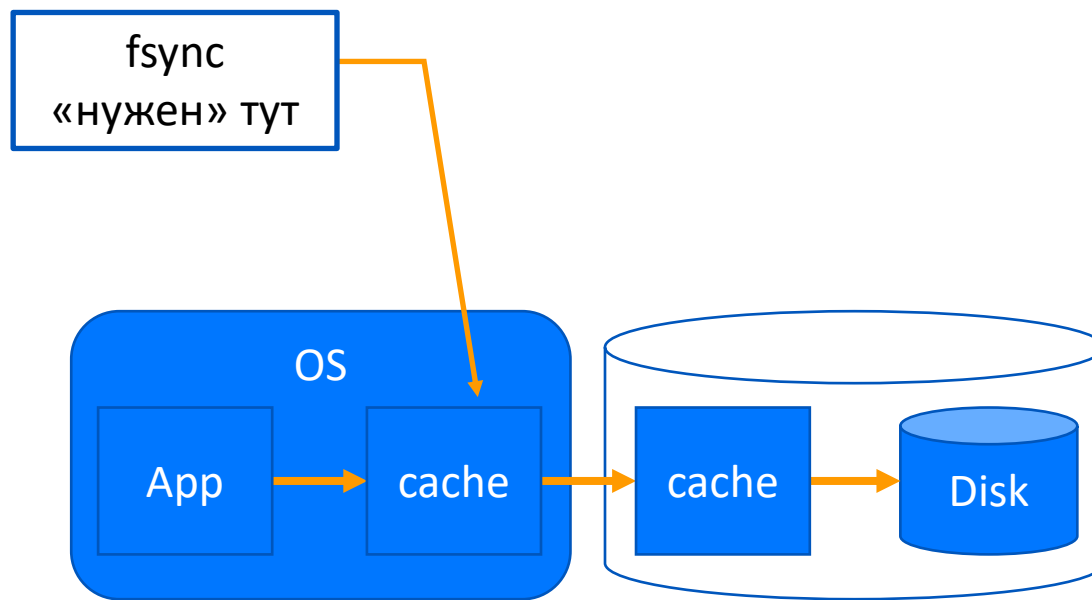
Для чего нужен fsync?



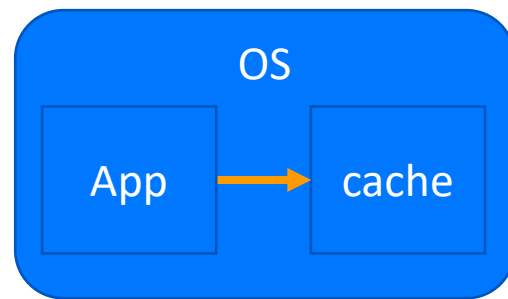
Для чего нужен fsync?



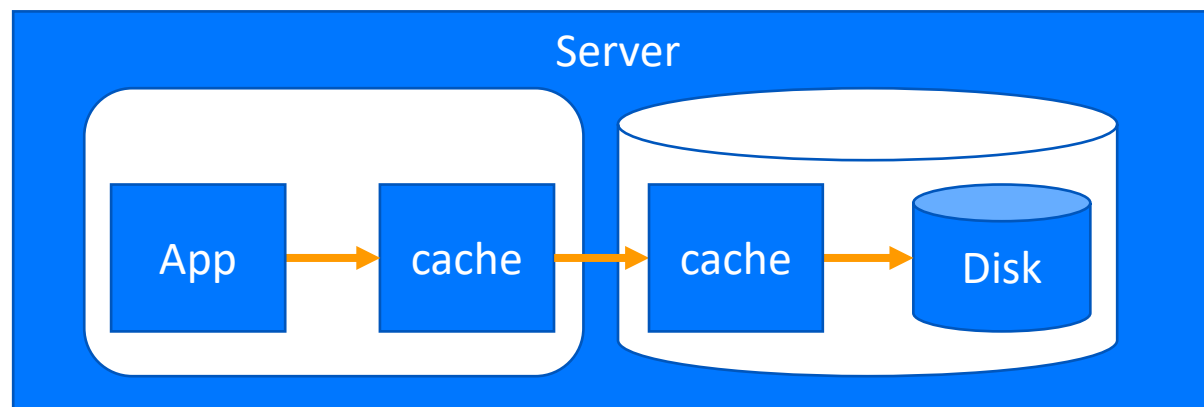
Для чего нужен fsync?



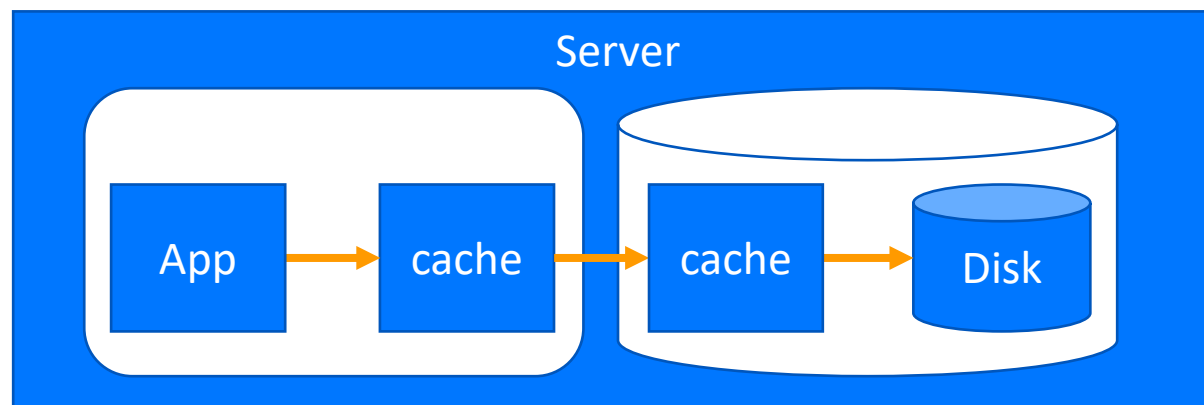
Для чего нужен fsync?



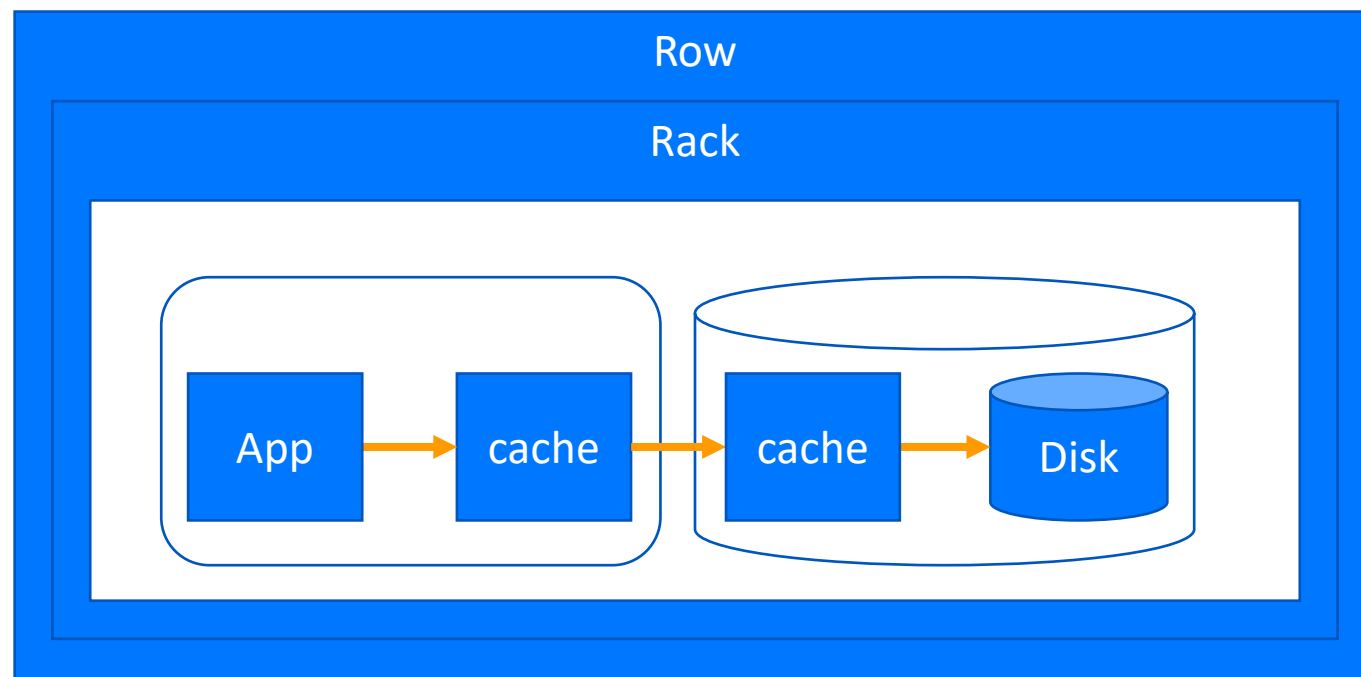
Для чего нужен fsync?



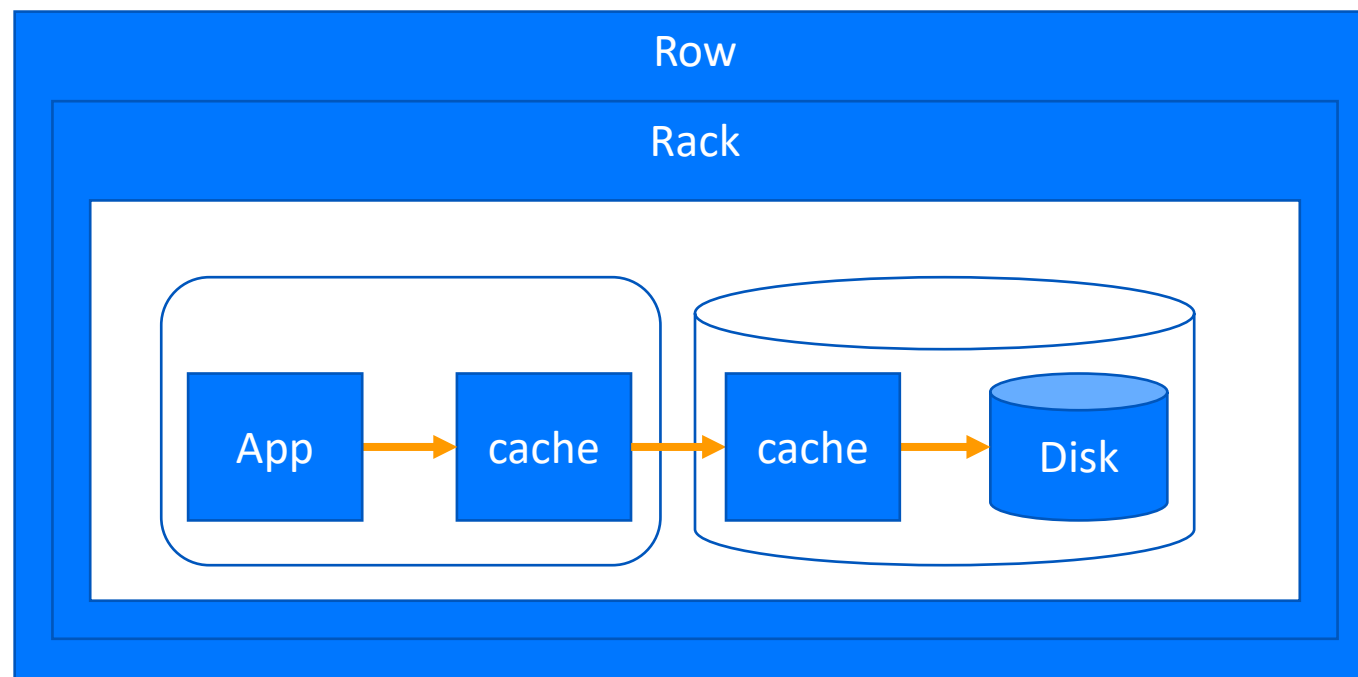
Для чего нужен fsync?



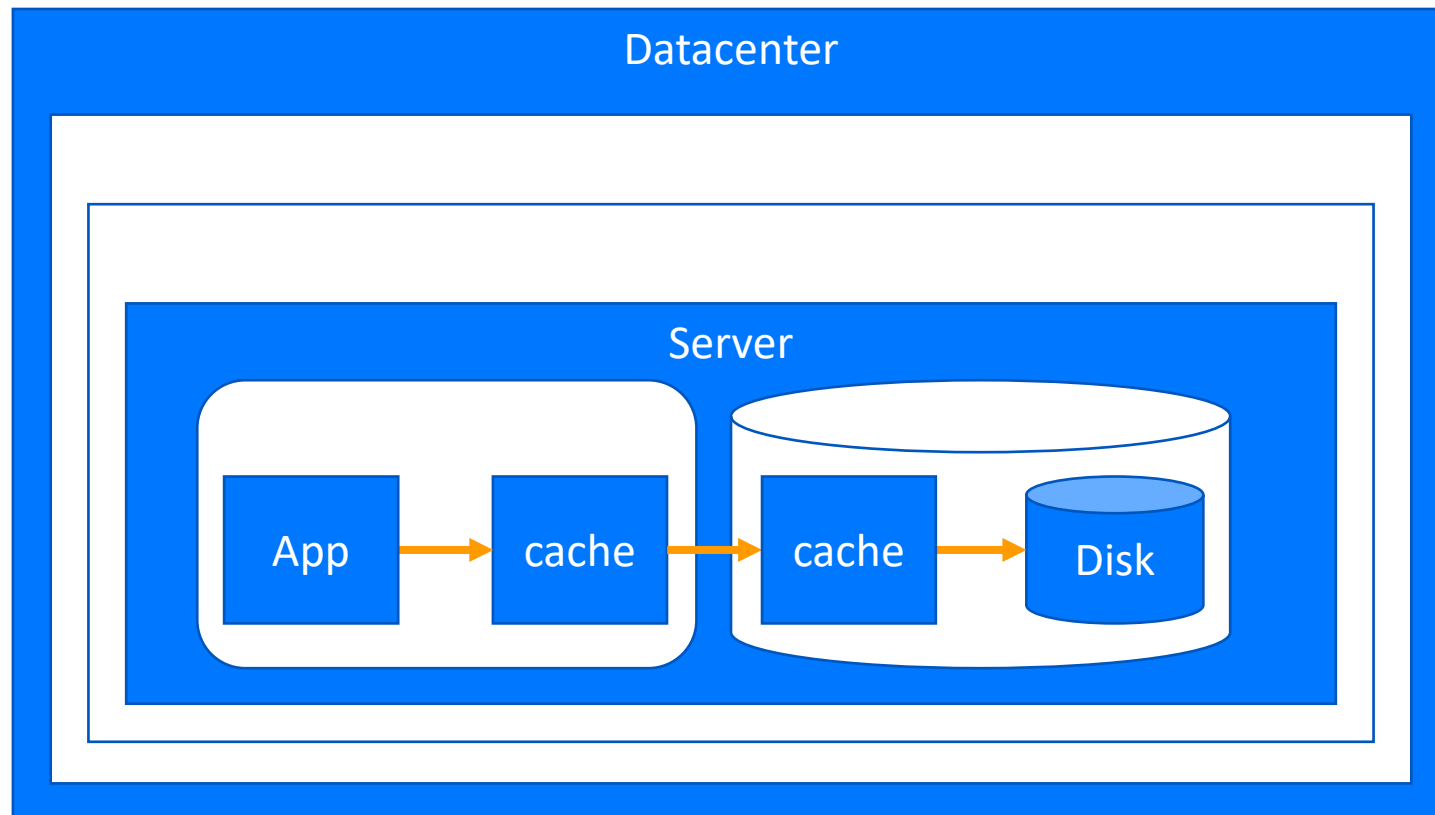
Для чего нужен fsync?



Для чего нужен fsync?



Для чего нужен fsync?



Датацентры горят



fsync защищает слишком слабо

Не защищает от отказа диска

Не защищает от отказа сервера

Не защищает от отказа датацентра

fsync защищает слишком слабо

Не защищает от отказа диска

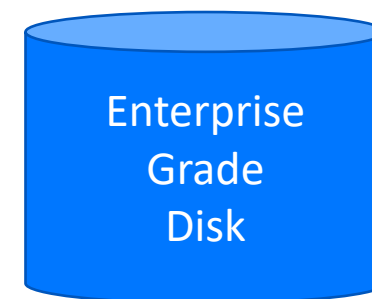
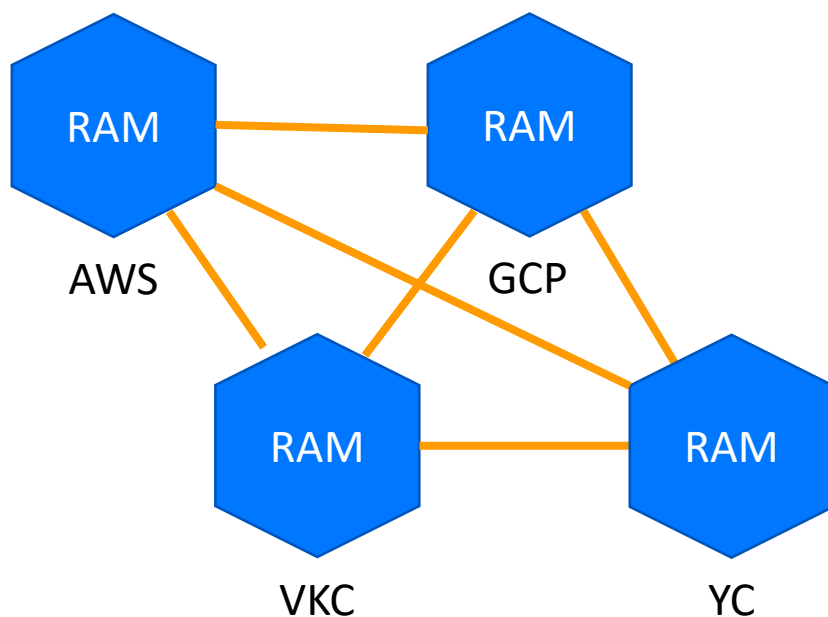
Не защищает от отказа сервера

Не защищает от отказа датацентра

Поэтому нужна репликация

Память против дисков

Распределённая и задублированная система, не сохраняющая на диск, будет надёжнее одного диска*



Вероятность одновременного отказа

$$P_n = P^n$$

- 1 — 50%

Вероятность одновременного отказа

$$P_n = P^n$$

- 1 — 50%
- 2 — 25%
- 3 — 12.5%
- 4 — 6.25%

Вероятность одновременного отказа

$$P_n = P^n$$

- 1 — 50%
- 2 — 25%
- 3 — 12.5%
- 4 — 6.25%
- 10 — 0.1%

Сохранность в памяти + репликация

Читаем из памяти

Пишем в память

Пишем на диск

Ждём подтверждения

Реплицируем

(Ждём подтверждения)

Сохранность в памяти + репликация

Потерять...
...все данные

RPO = «начало»

Потерять...
... «последние» данные

RPO = «последние»

Сохранность в памяти + репликация

Потерять...
...все данные

RPO = «начало»

Потерять...
... «последние» данные

RPO = «последние»

* RPO — Recovery Point Objective

Сохранность в памяти + репликация

Потерять...
...все данные

RPO = «начало»

Потерять...
... «последние» данные

RPO = лаг репликации

* RPO — Recovery Point Objective

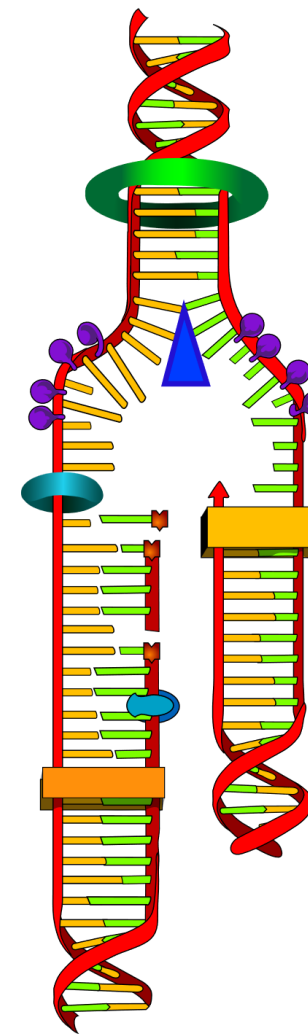
04

Репликация

Какие виды репликации бывают и почему так много?

Что такое репликация?

Механизм создания точной копии



Зачем нужна репликация?

Резервирование:

- Одно-нодовые системы ненадежны (диски, питание)
- Одного ДЦ недостаточно (пожар)

Зачем нужна репликация?

Резервирование:

- Одно-нодовые системы ненадежны (диски, питание)
- Одного ДЦ недостаточно (пожар)

Доступность:

- Время восстановления после аварии (RTO)

Зачем нужна репликация?

Резервирование:

- Одно-нодовые системы ненадежны (диски, питание)
- Одного ДЦ недостаточно (пожар)

Доступность:

- Время восстановления после аварии (RTO)

Производительность:

- Железо узкое место (1 thread, RAM)

Какая бывает репликация

Синхронная — надёжность

Асинхронная — скорость

Синхронная и асинхронная

Синхронная

подтверждается после фиксации
на реплике

надежная (durable)

медленная

выше шанс отказа

RPO = 0

Асинхронная

подтверждается после записи на
одну ноду

можно потерять часть

«быстрая»

отказ реплики не важен

RPO > 0 (\approx Lag)

* RPO — Recovery Point Objective

Какая бывает репликация

Однонаправленная — консистентность

Двунаправленная — доступность

Master-master vs Master-replica

MM

сложно использовать

RTO = 0

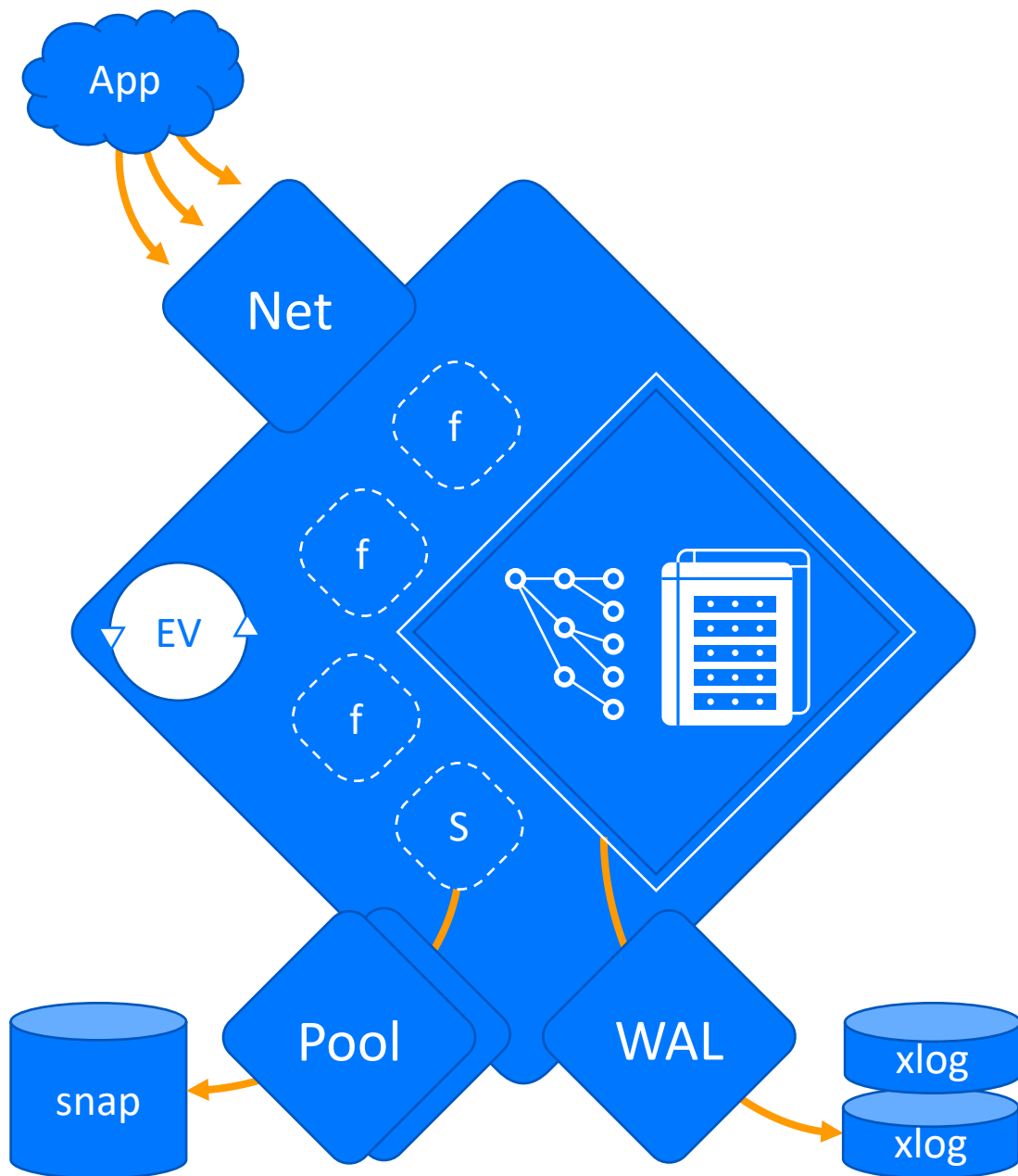
MR

просто использовать

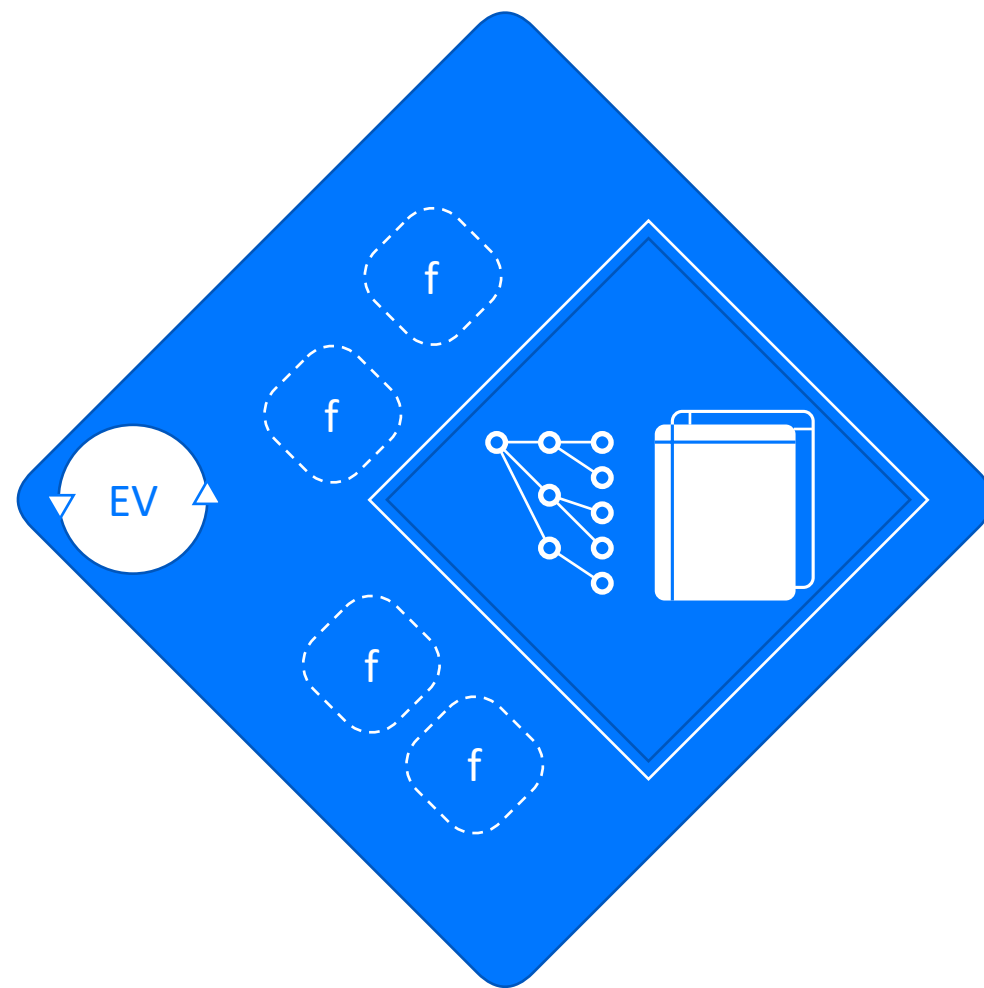
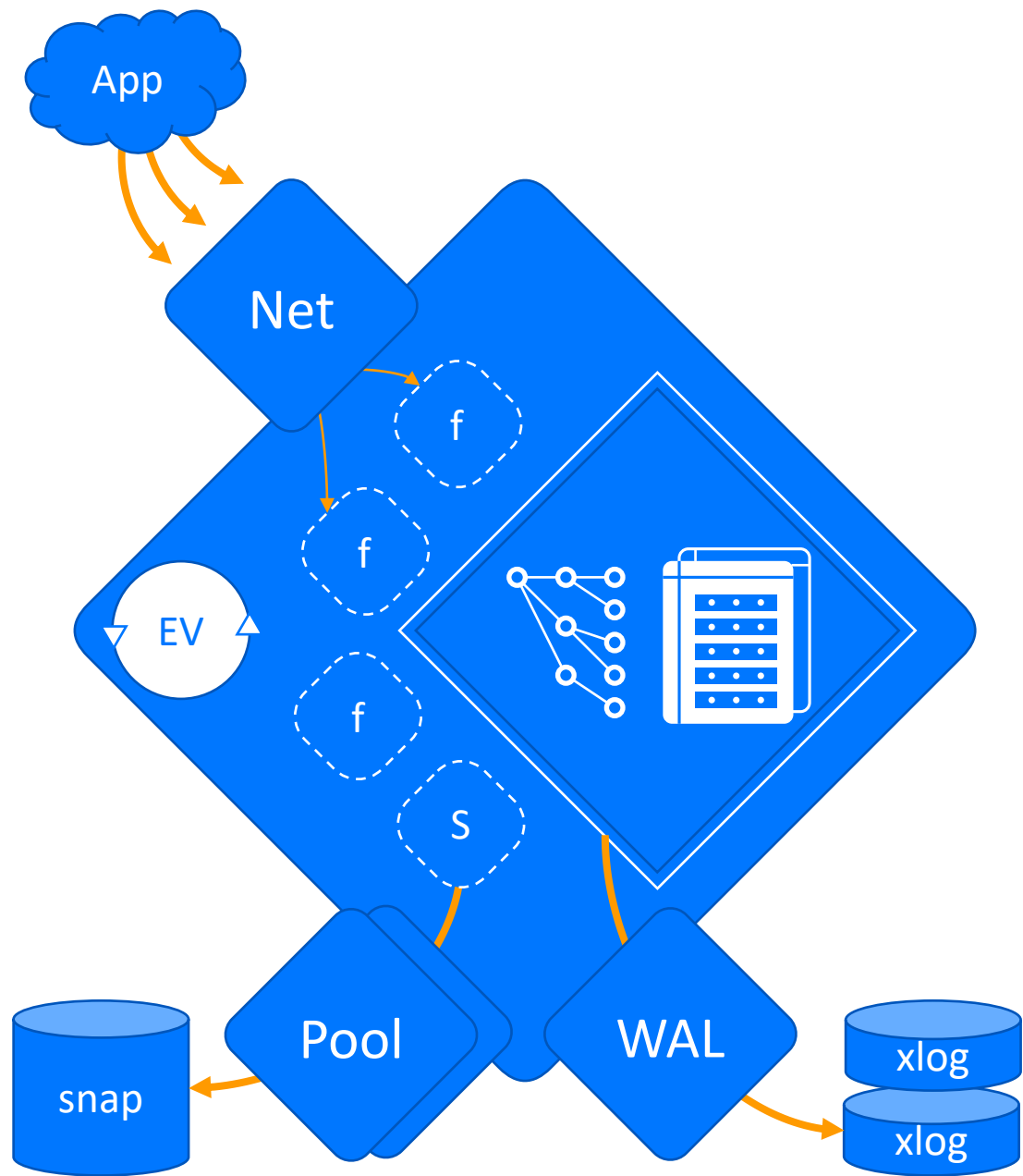
RTO > 0 (\approx Failover time)

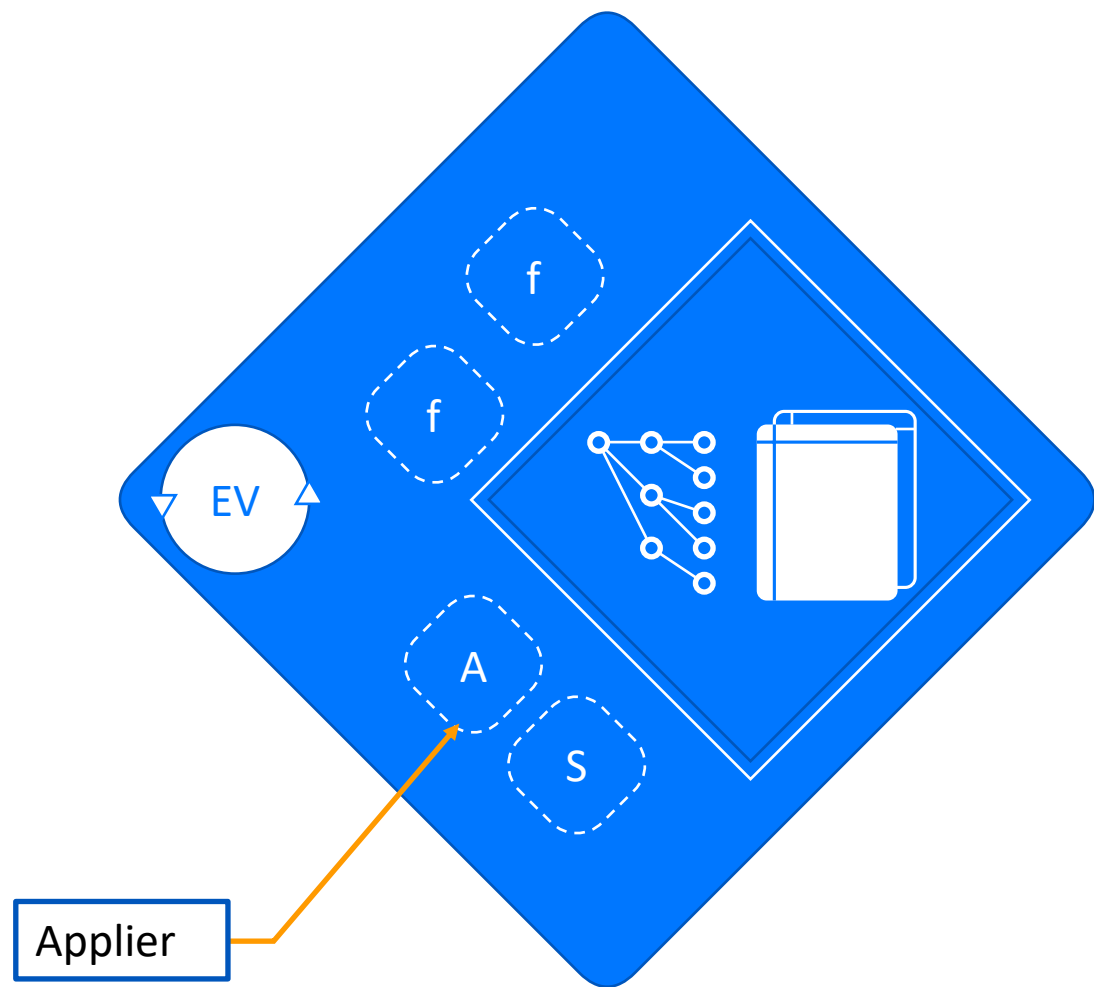
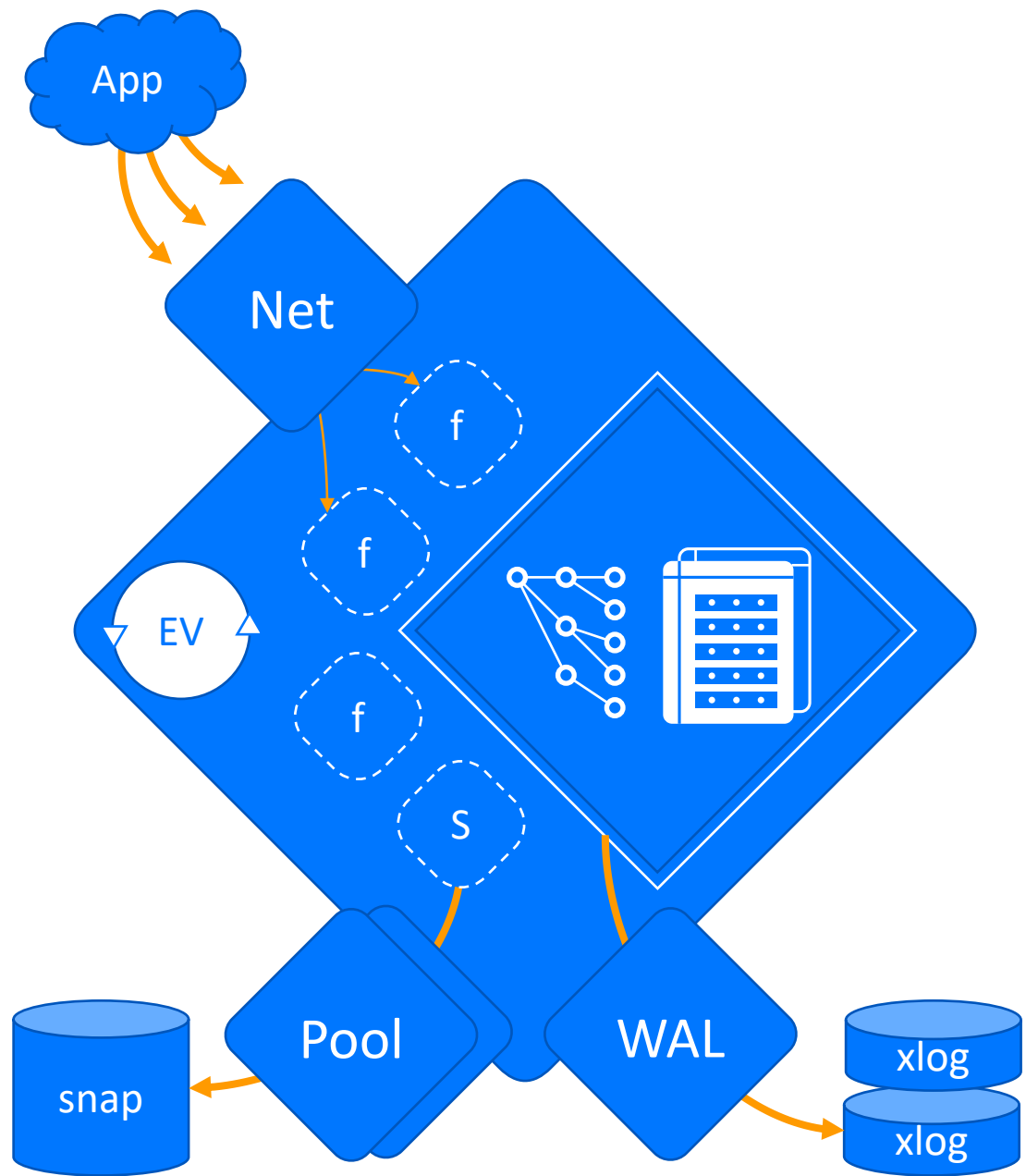
* RTO — Recovery Time Objective

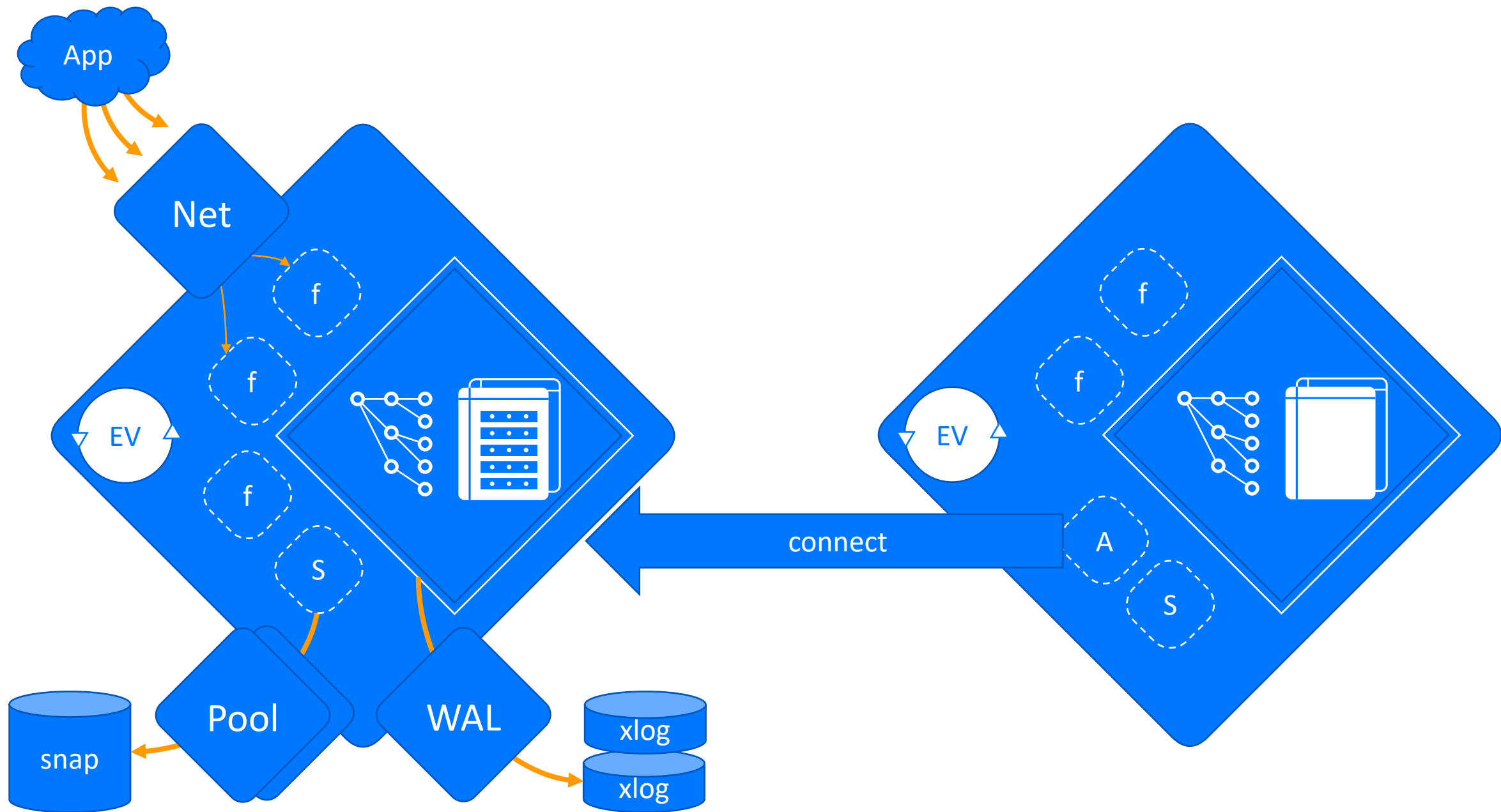
Устройство репликации в Tarantool

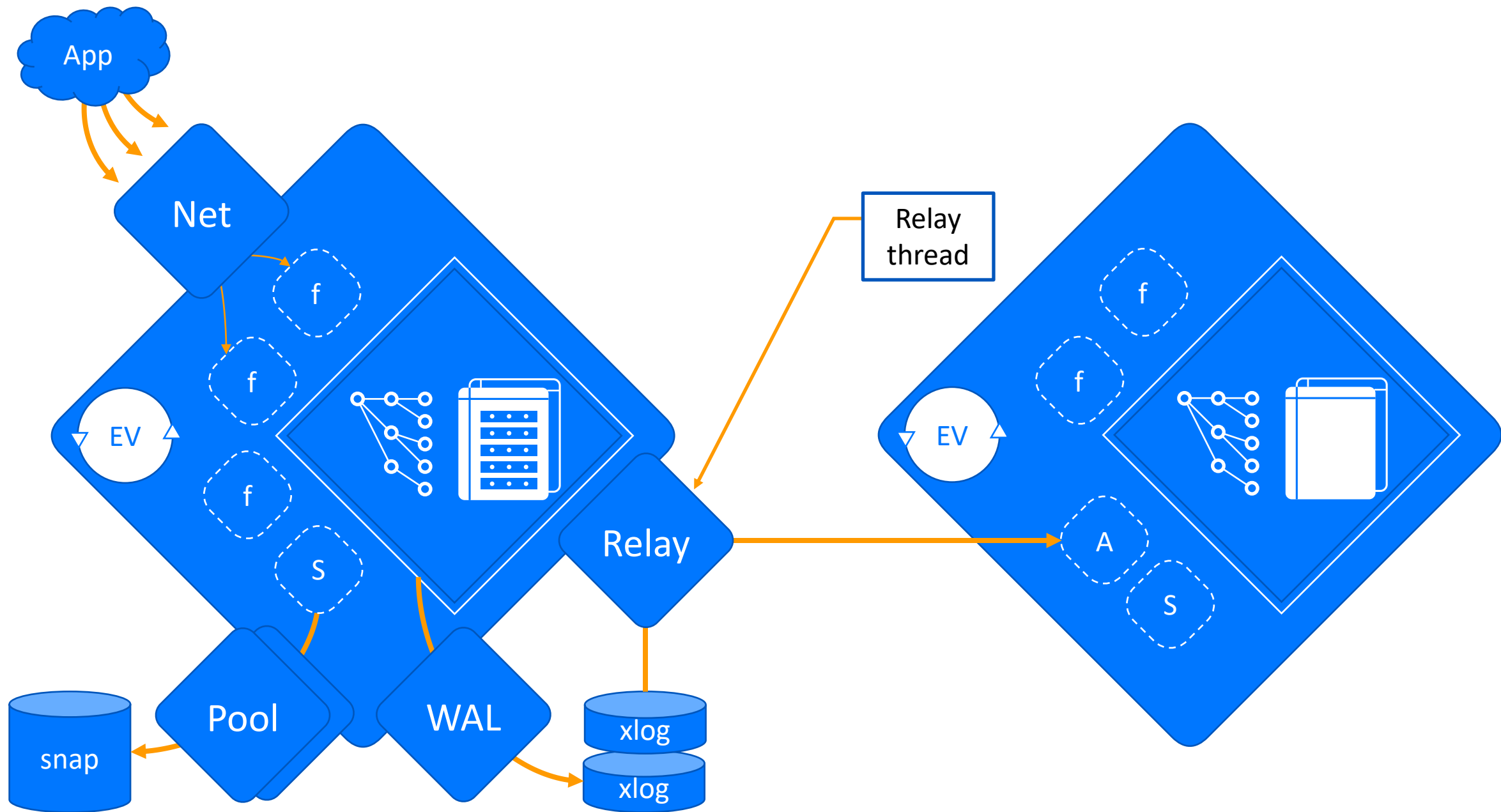


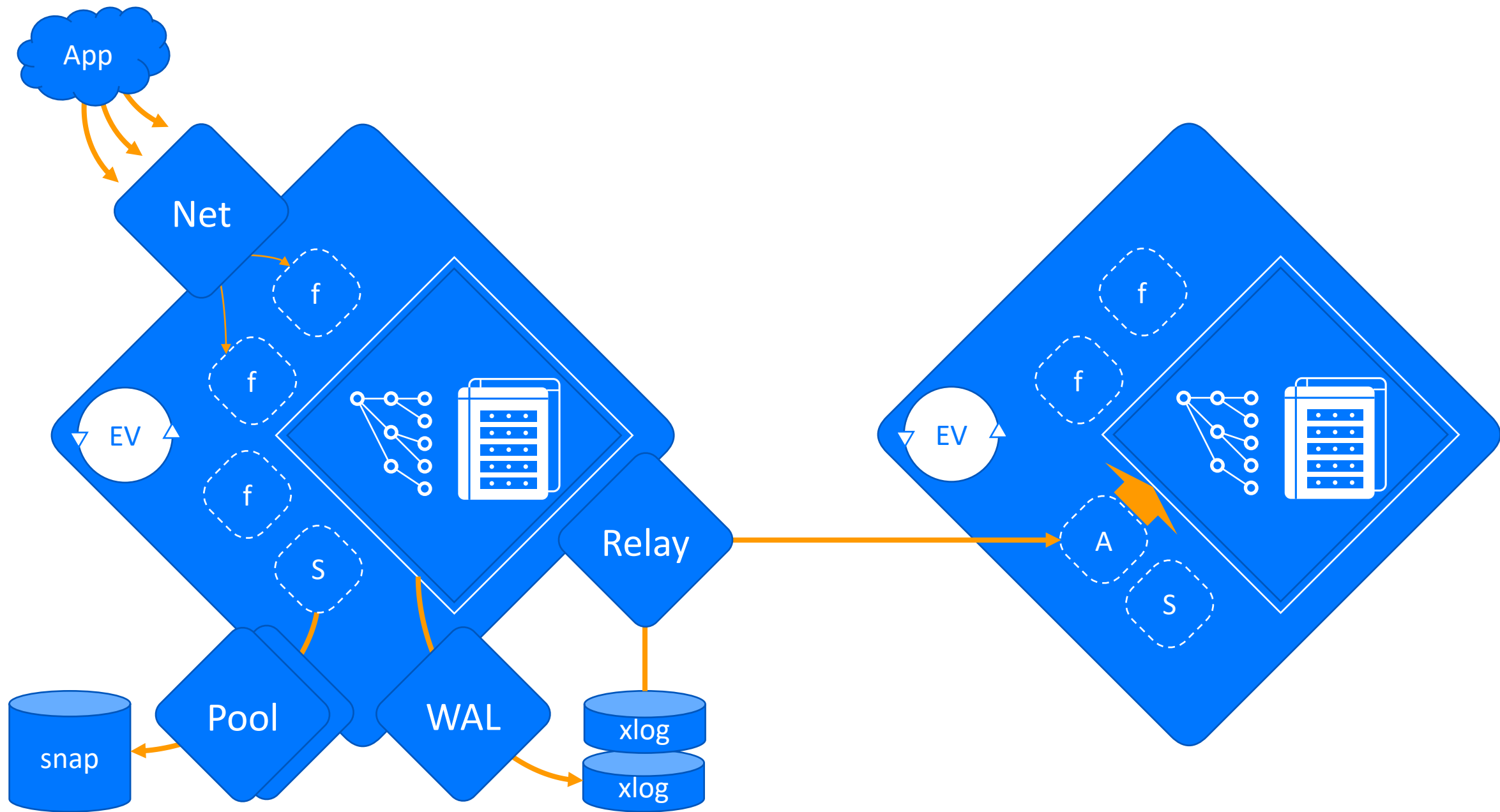
- Данные в памяти
- Монопольный доступ
- Событийный цикл
- Кооперативная многозадачность
- Отдельные потоки для сети
- Изменения во Write-Ahead-Log
- Консистентный снапшот
- Отдельные потоки для диска
- **WAL реплицируется**

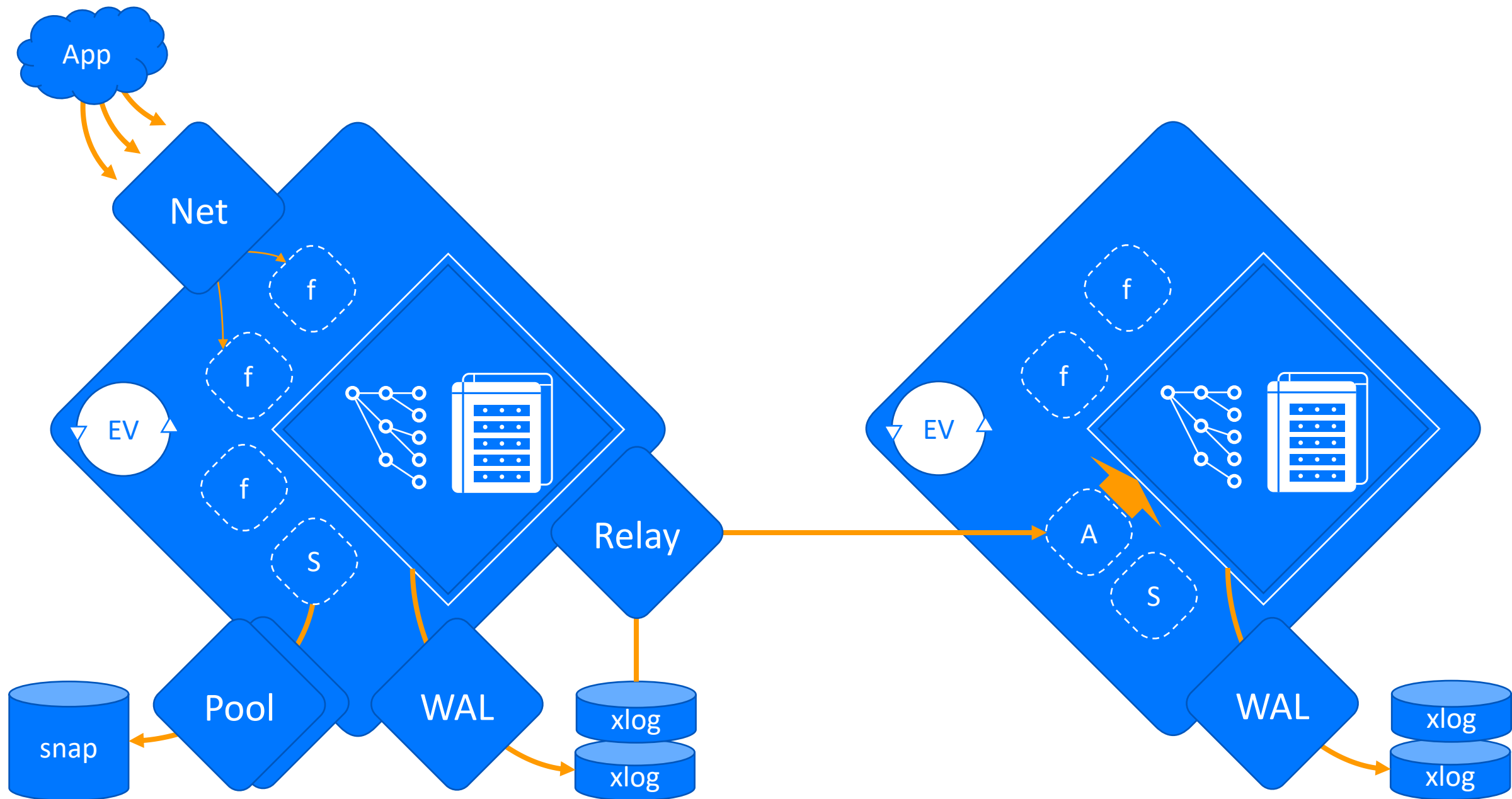


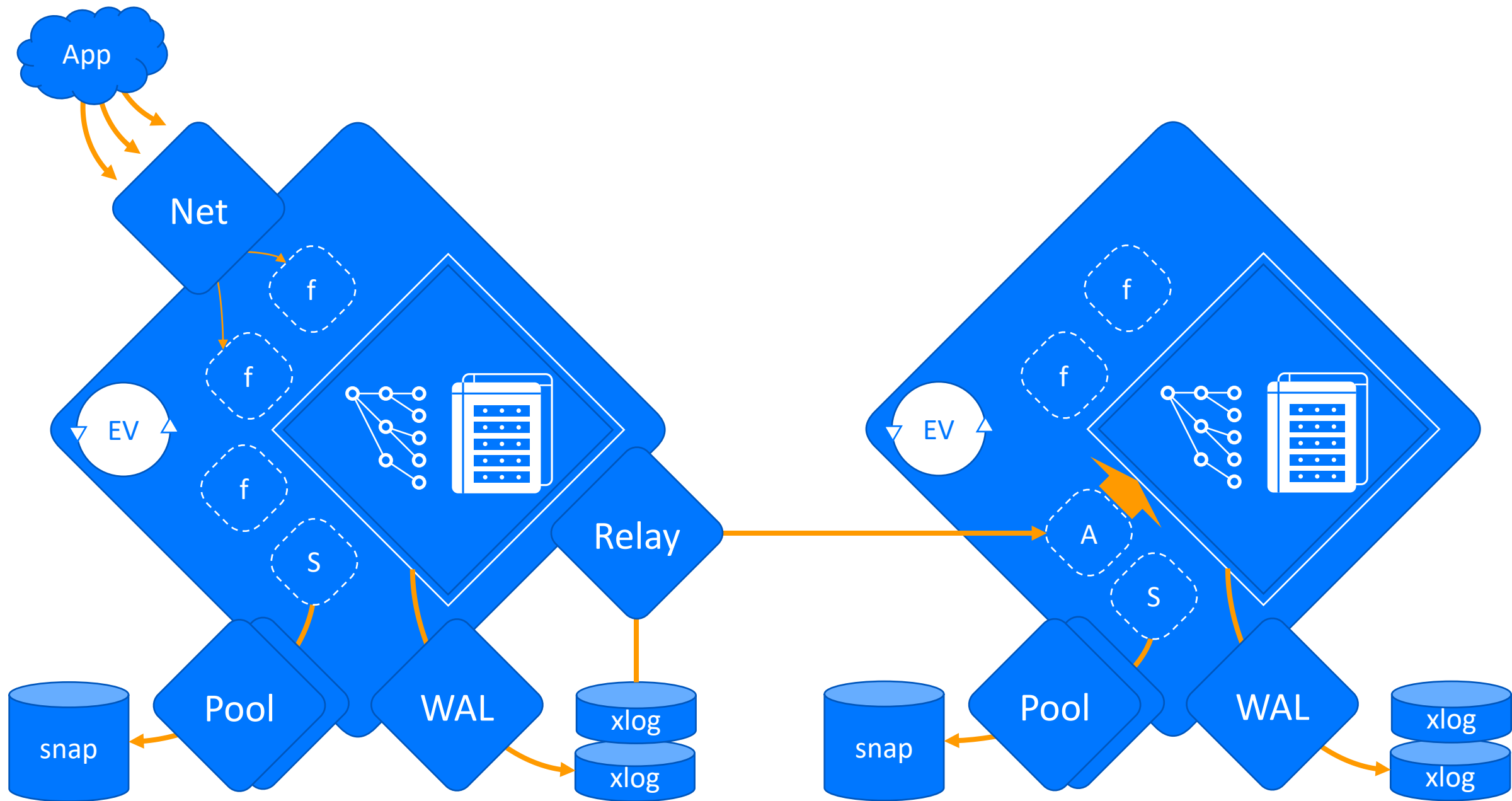


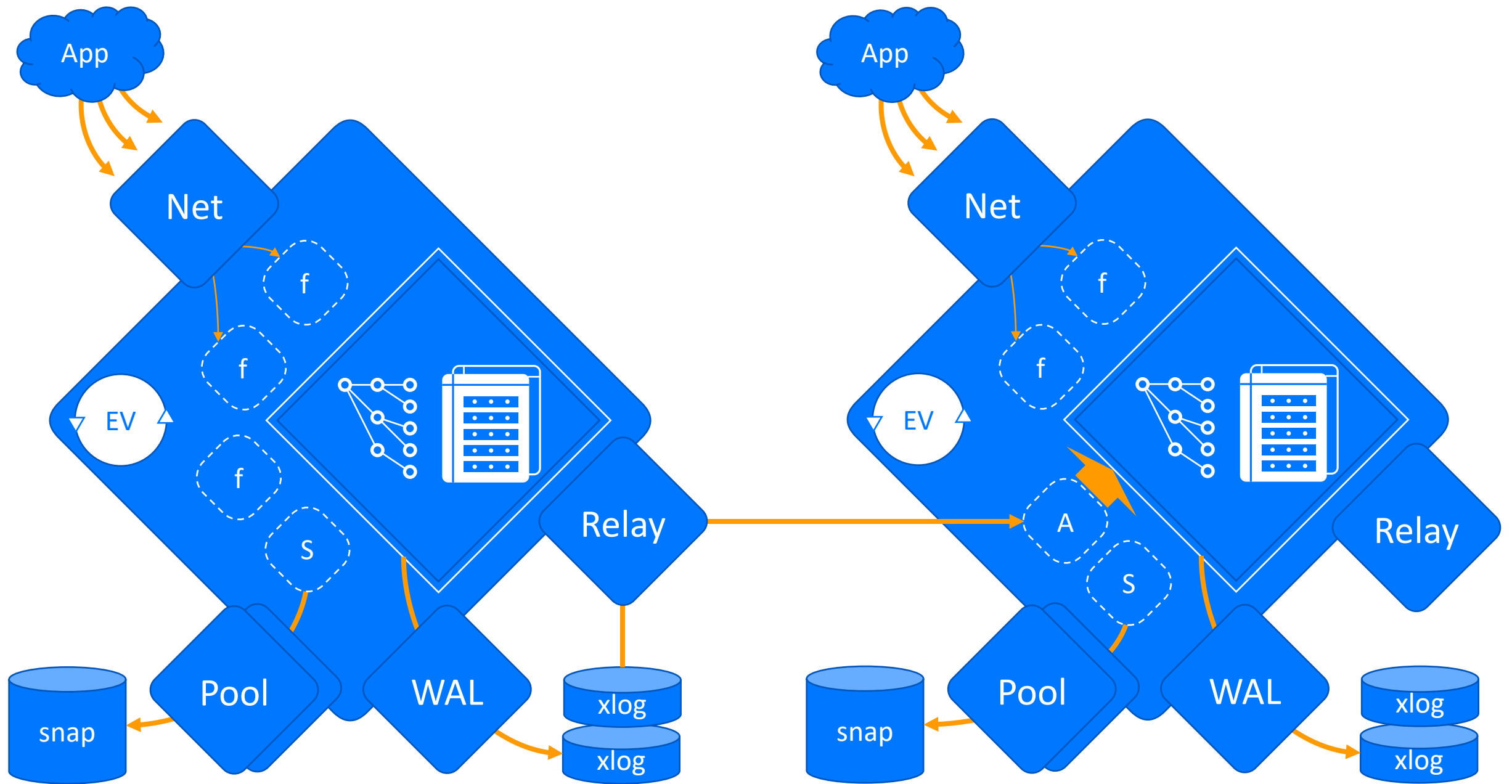


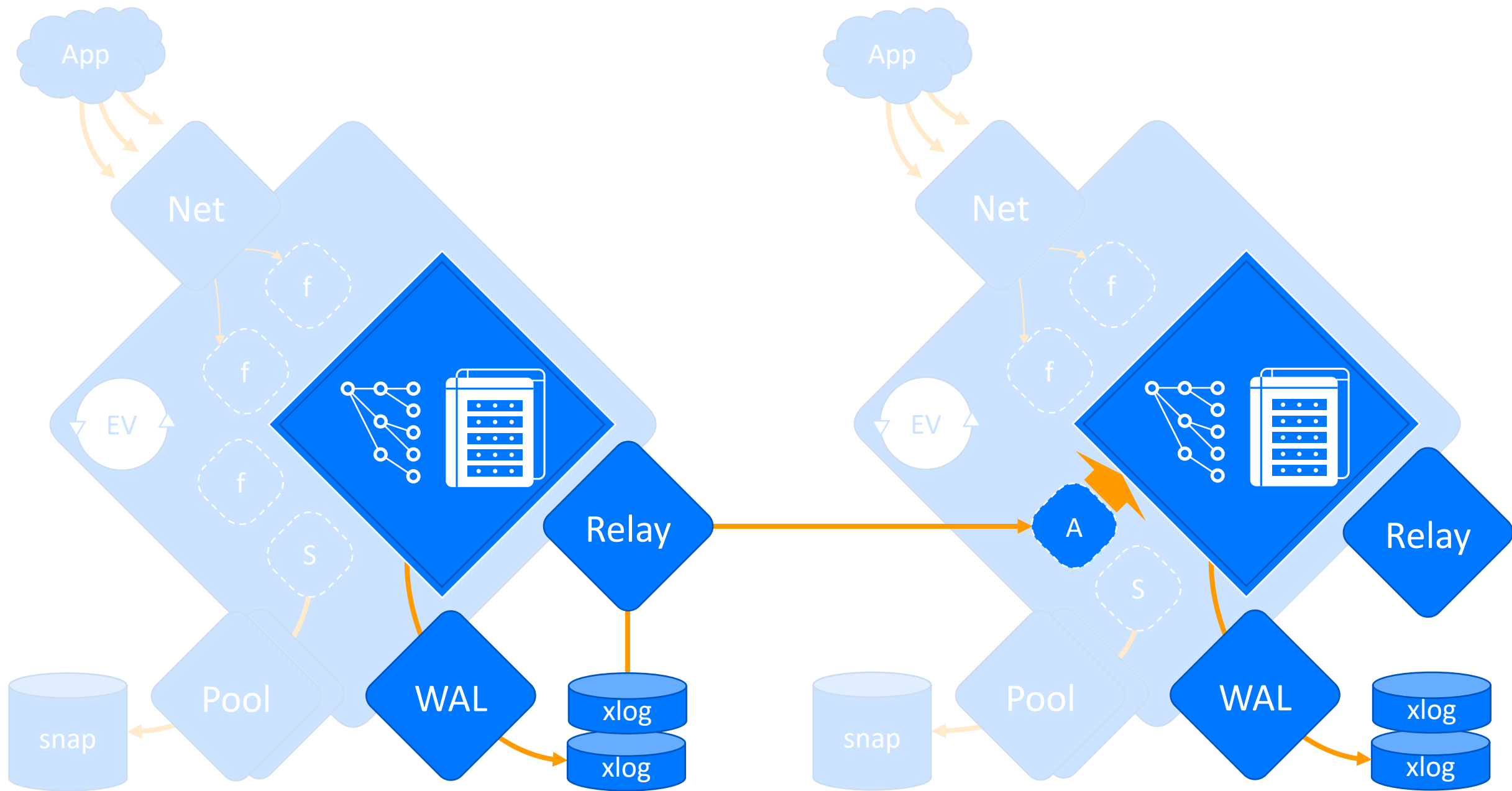








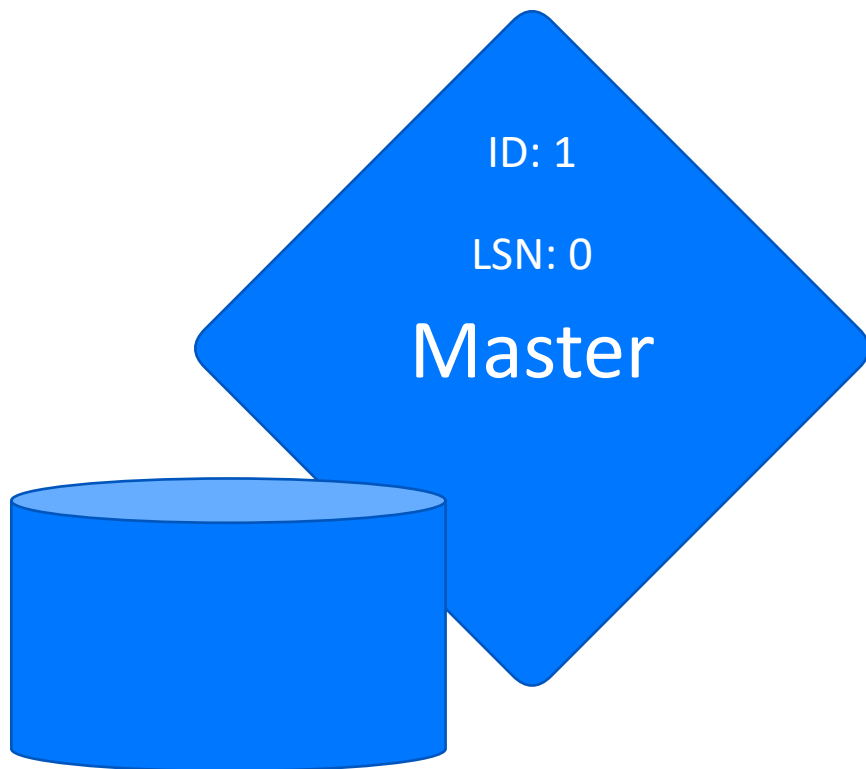




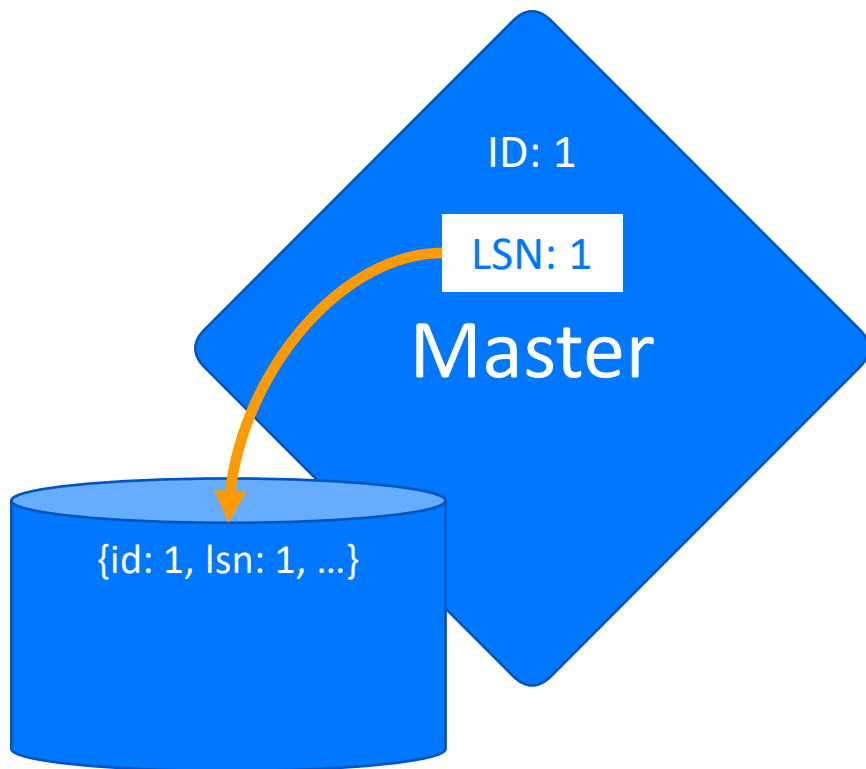
Log Sequence Number (LSN)

- Монотонная последовательность
- Определяет взаимный порядок транзакций
- Вместе с узлом — идентификатор транзакции

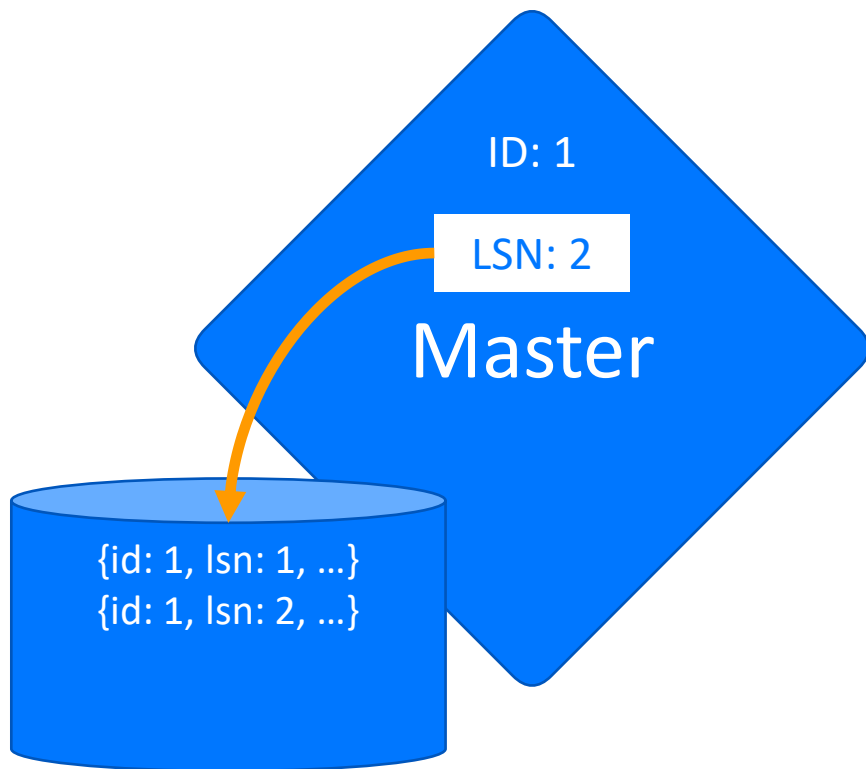
Репликация в Tarantool



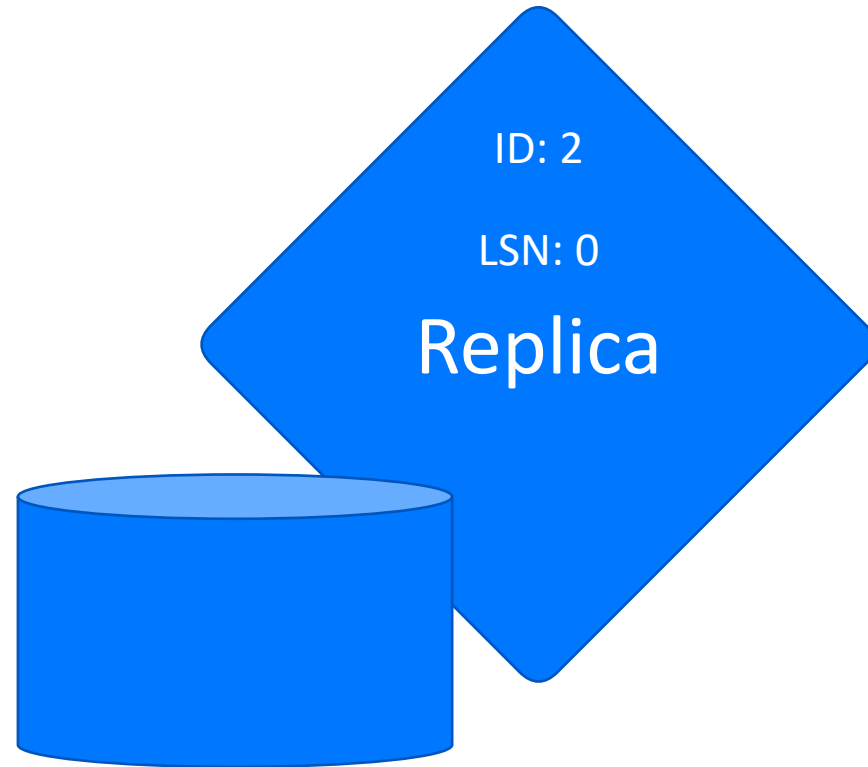
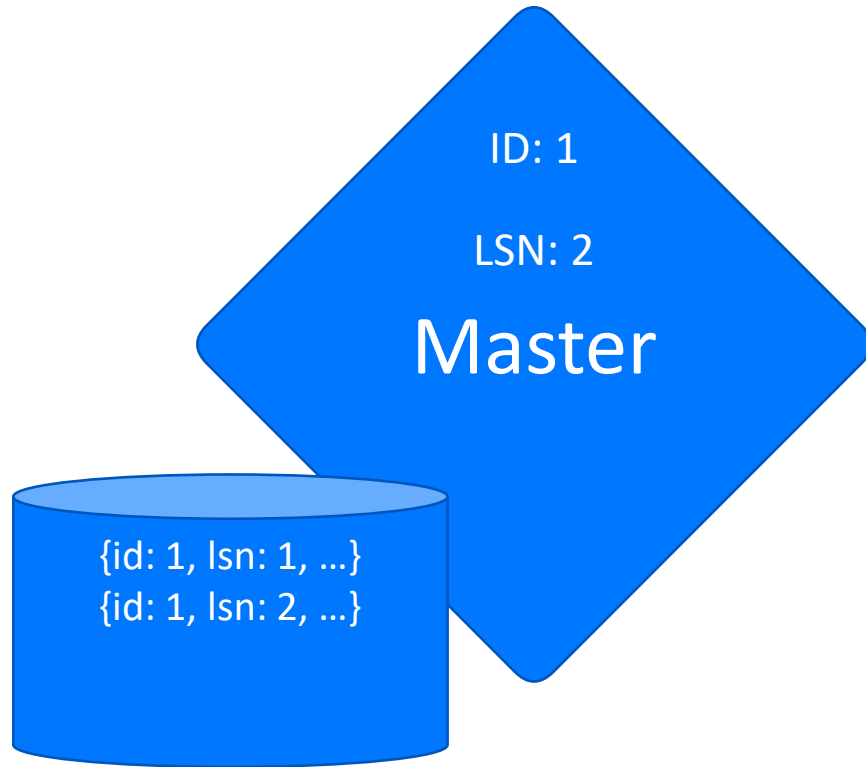
Репликация в Tarantool



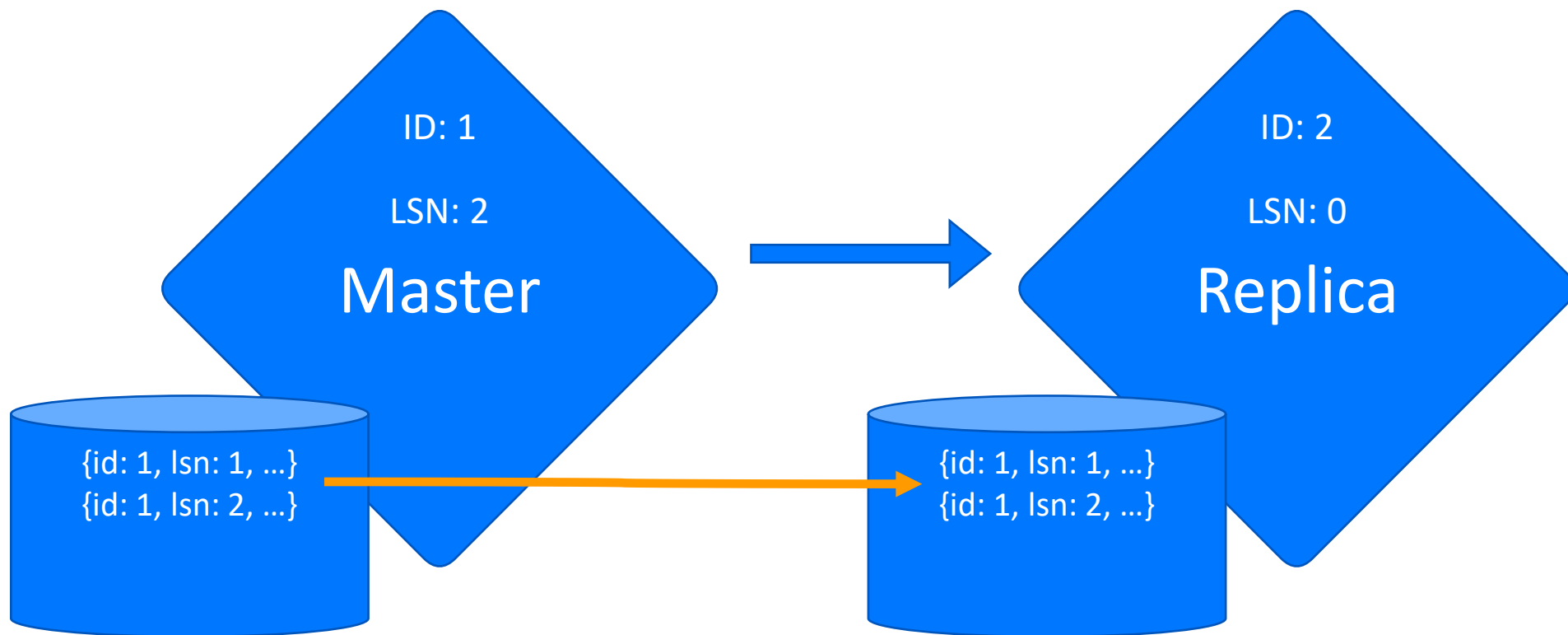
Репликация в Tarantool



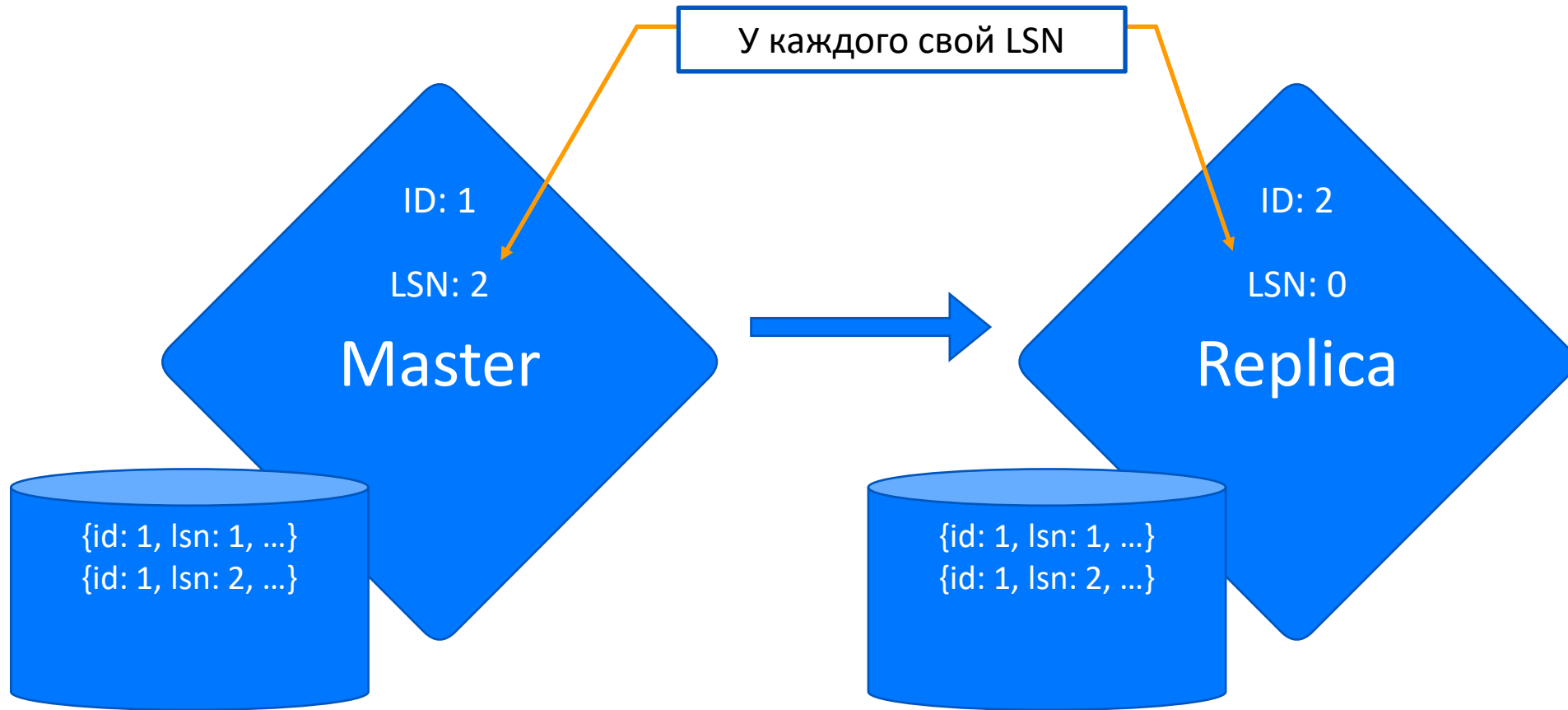
Репликация в Tarantool



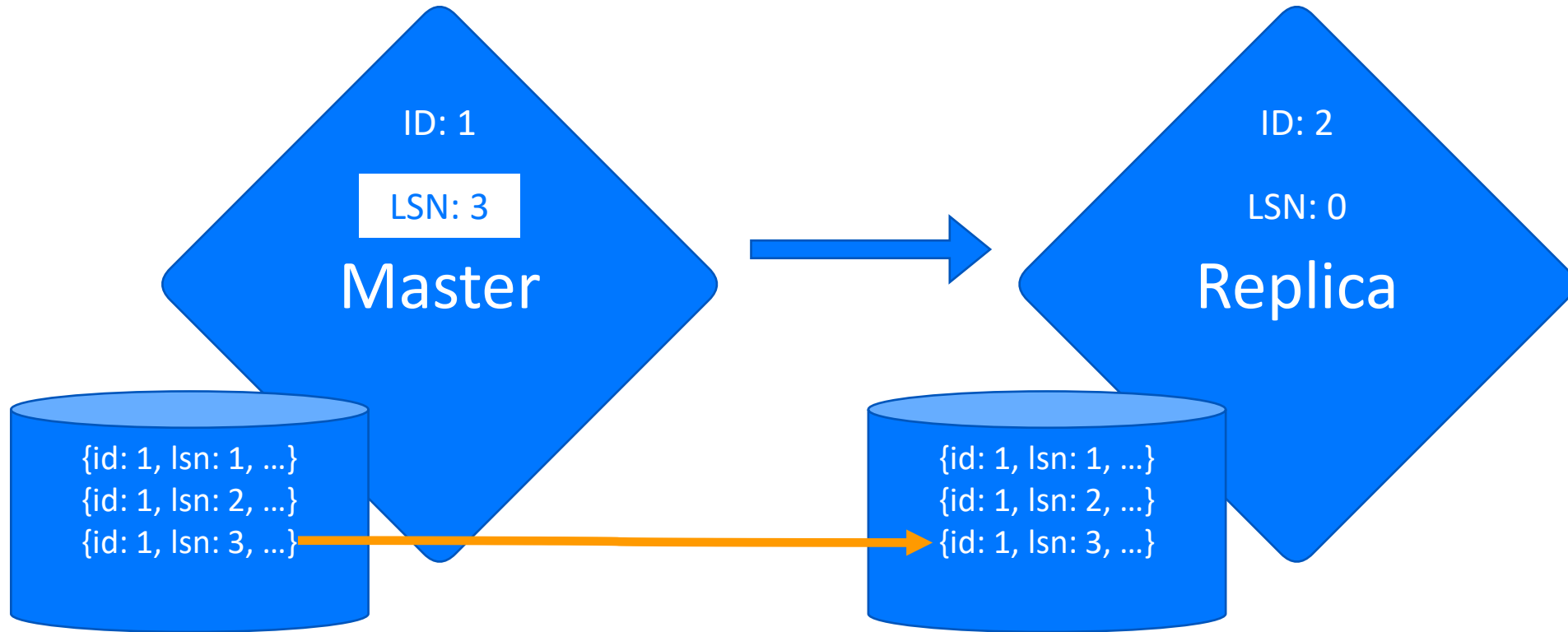
Репликация в Tarantool



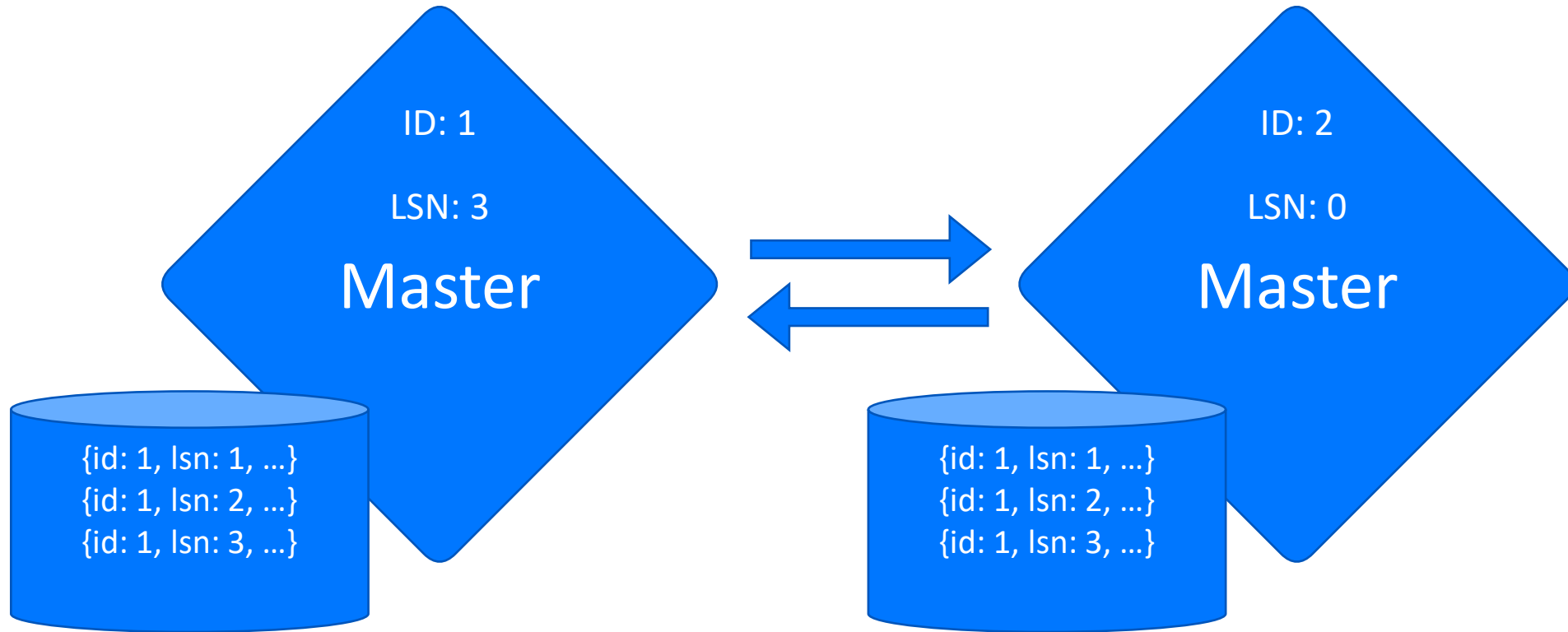
Репликация в Tarantool



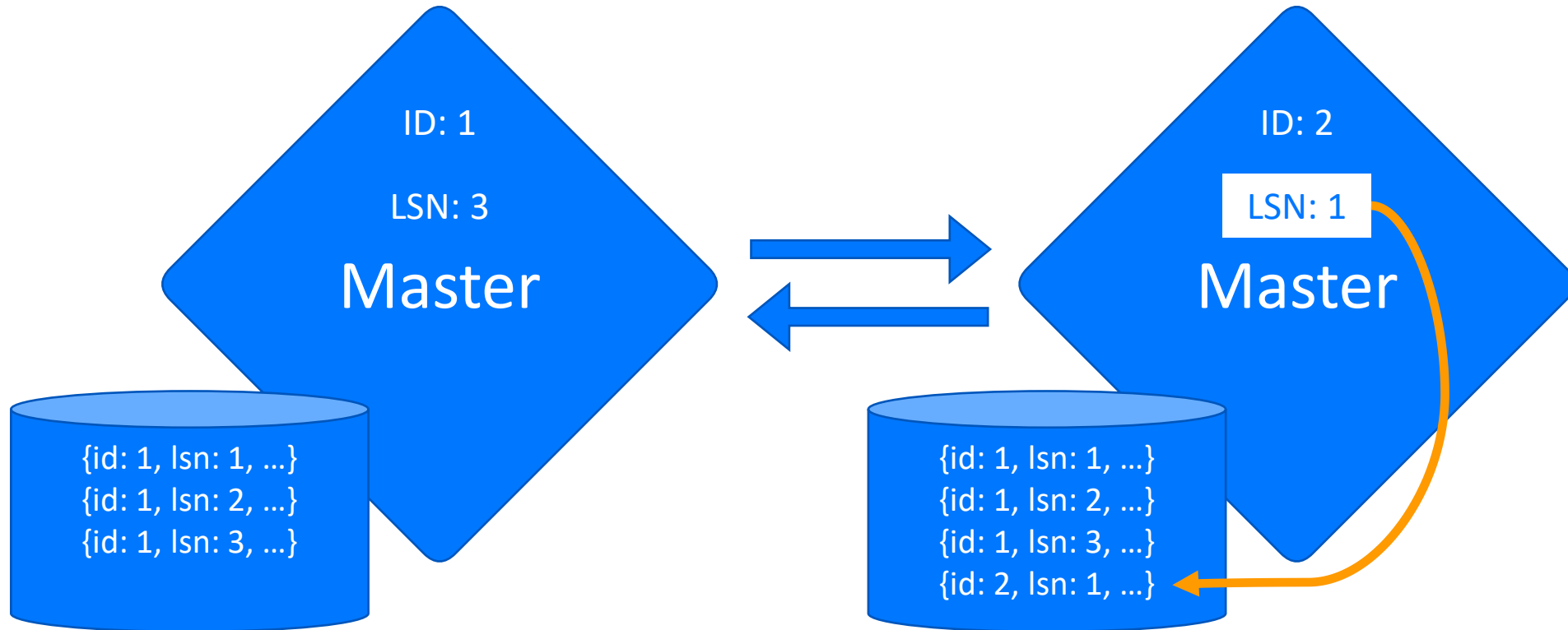
Репликация в Tarantool



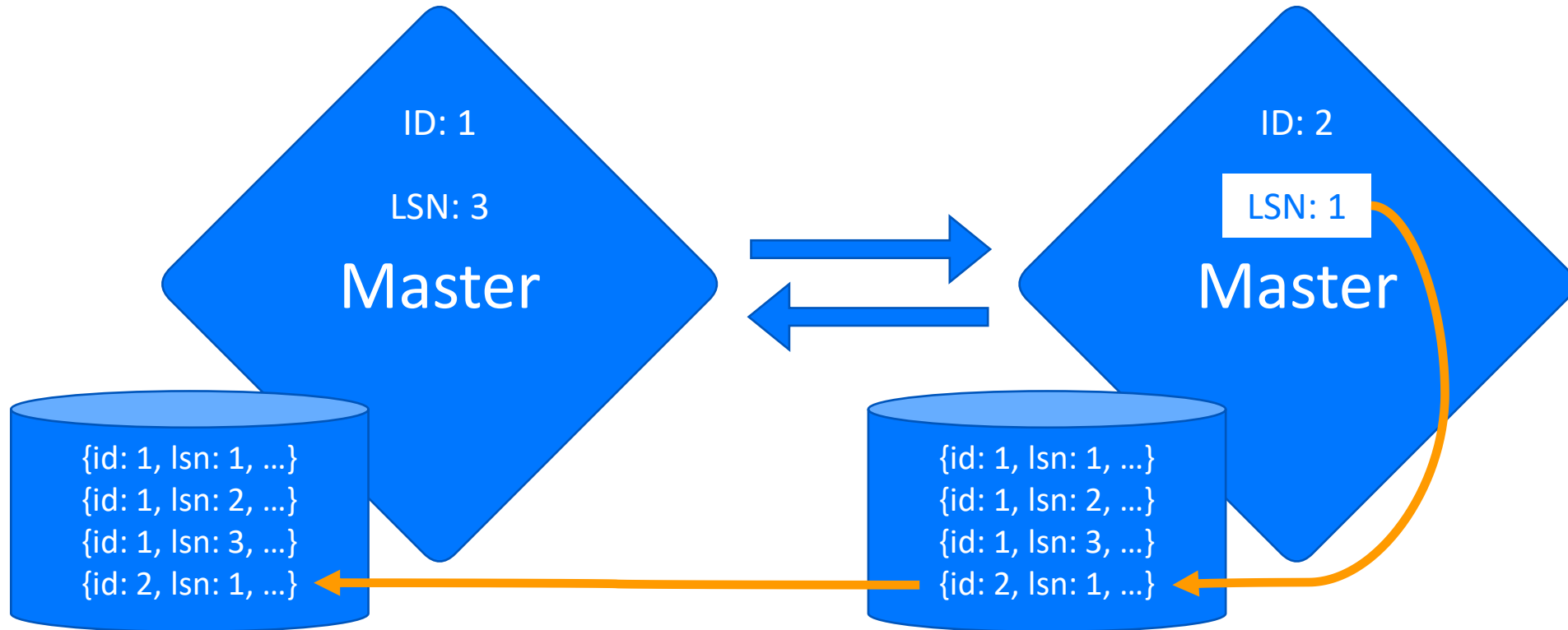
Репликация в Tarantool



Репликация в Tarantool



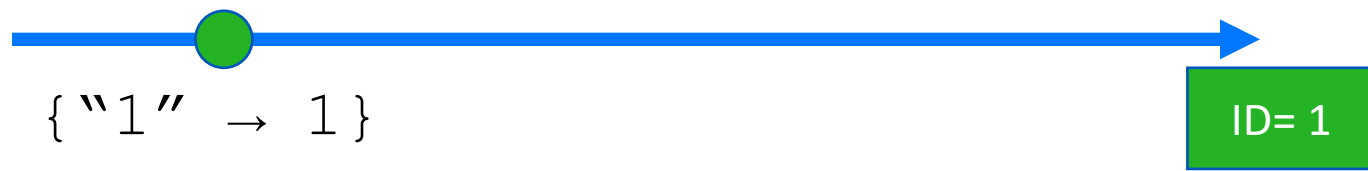
Репликация в Tarantool



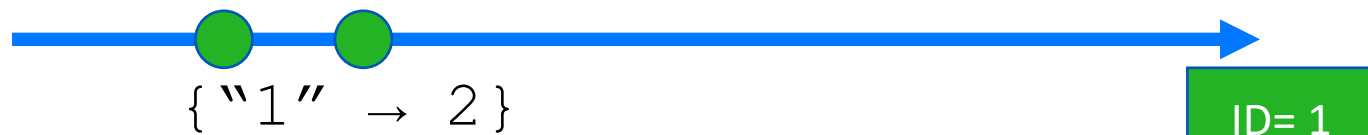
Векторные часы (vclock)

Алгоритм и структура
для упорядочивания событий
в распределённой системе

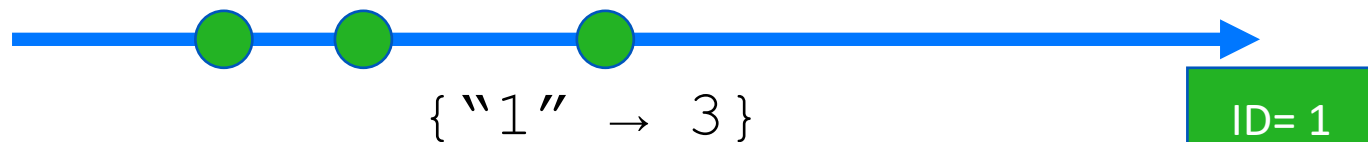
Векторные часы (vclock)



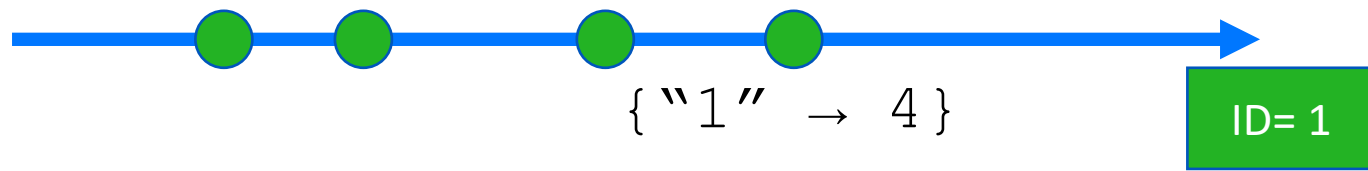
Векторные часы (vclock)



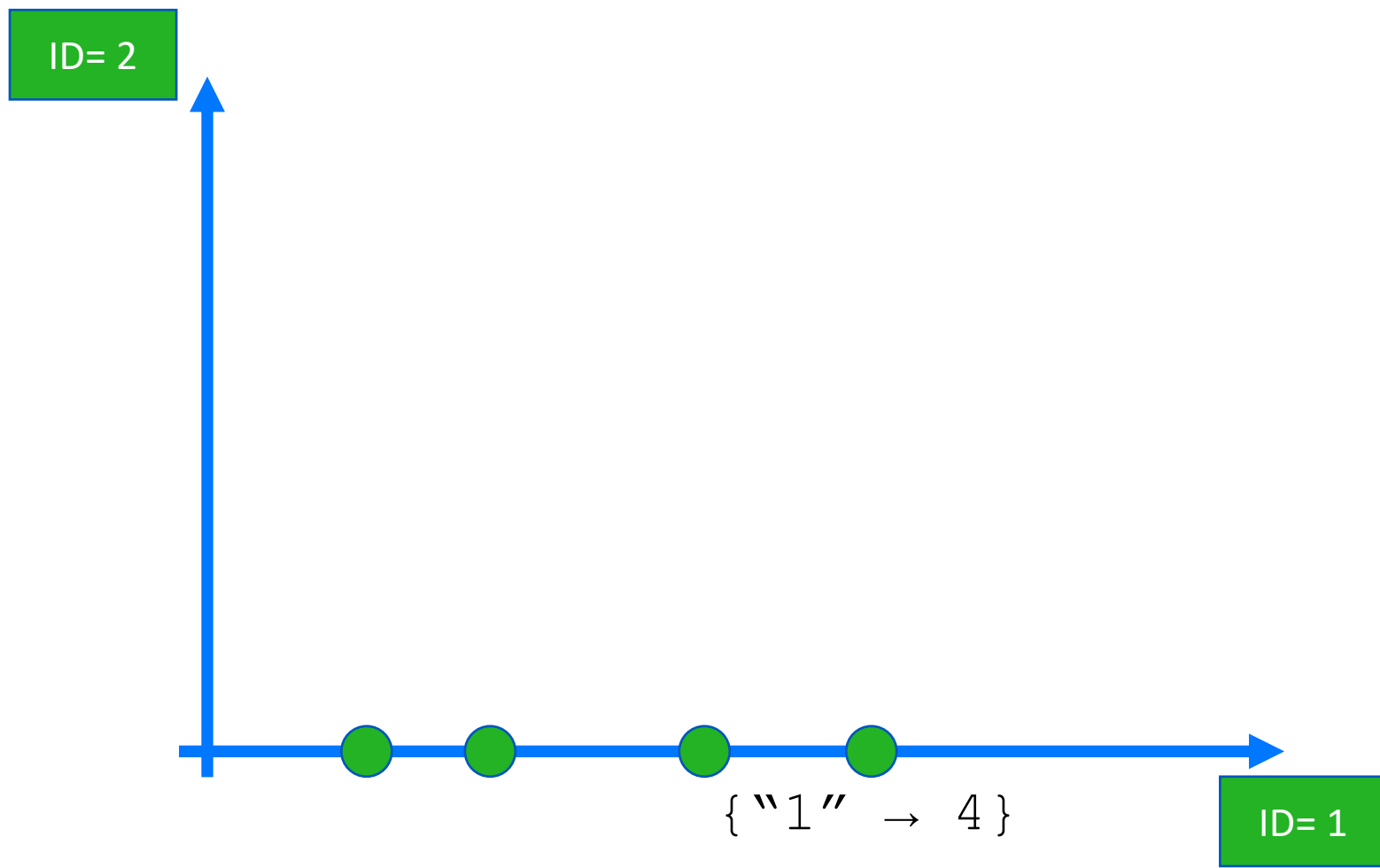
Векторные часы (vclock)



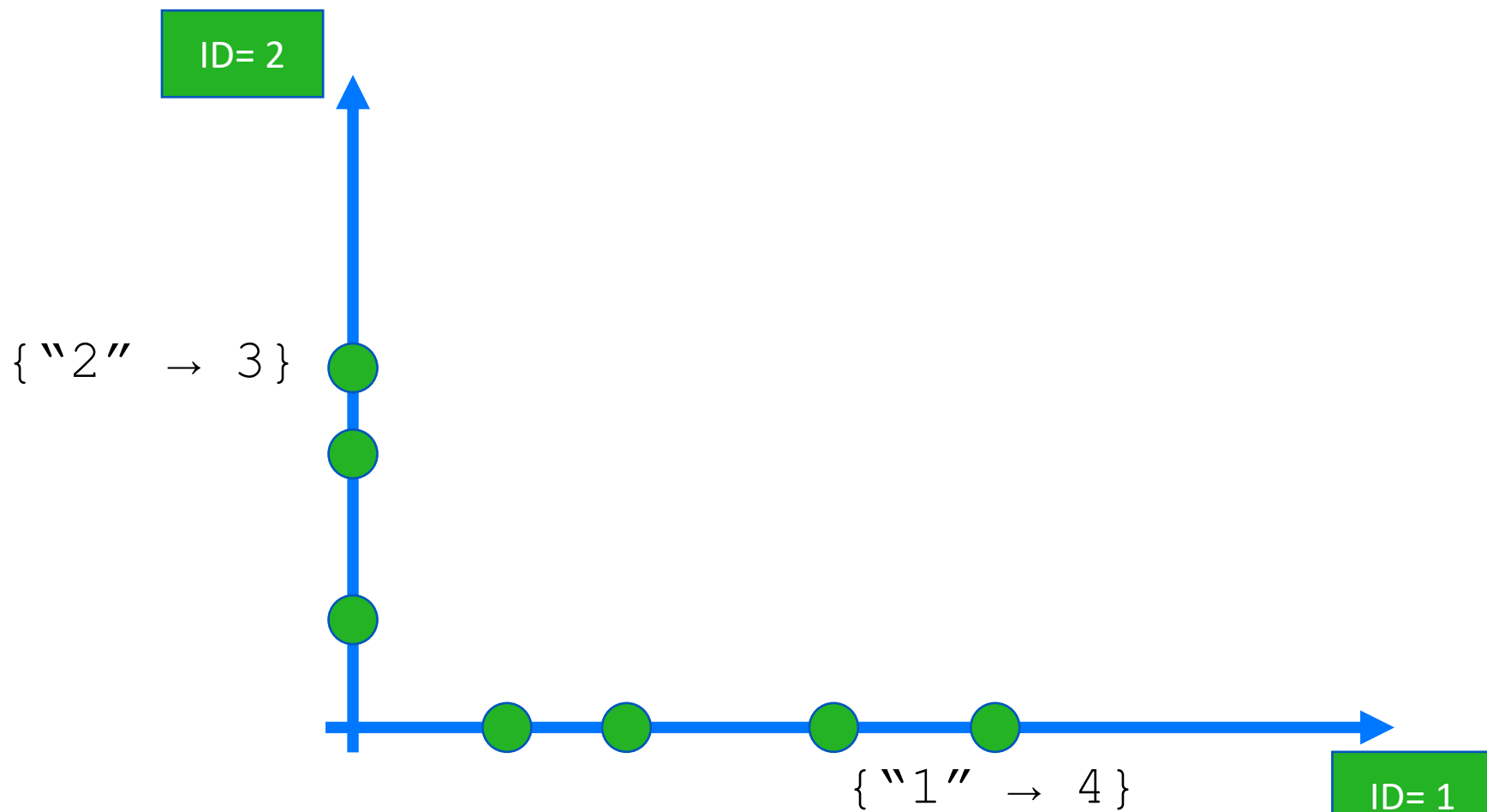
Векторные часы (vclock)



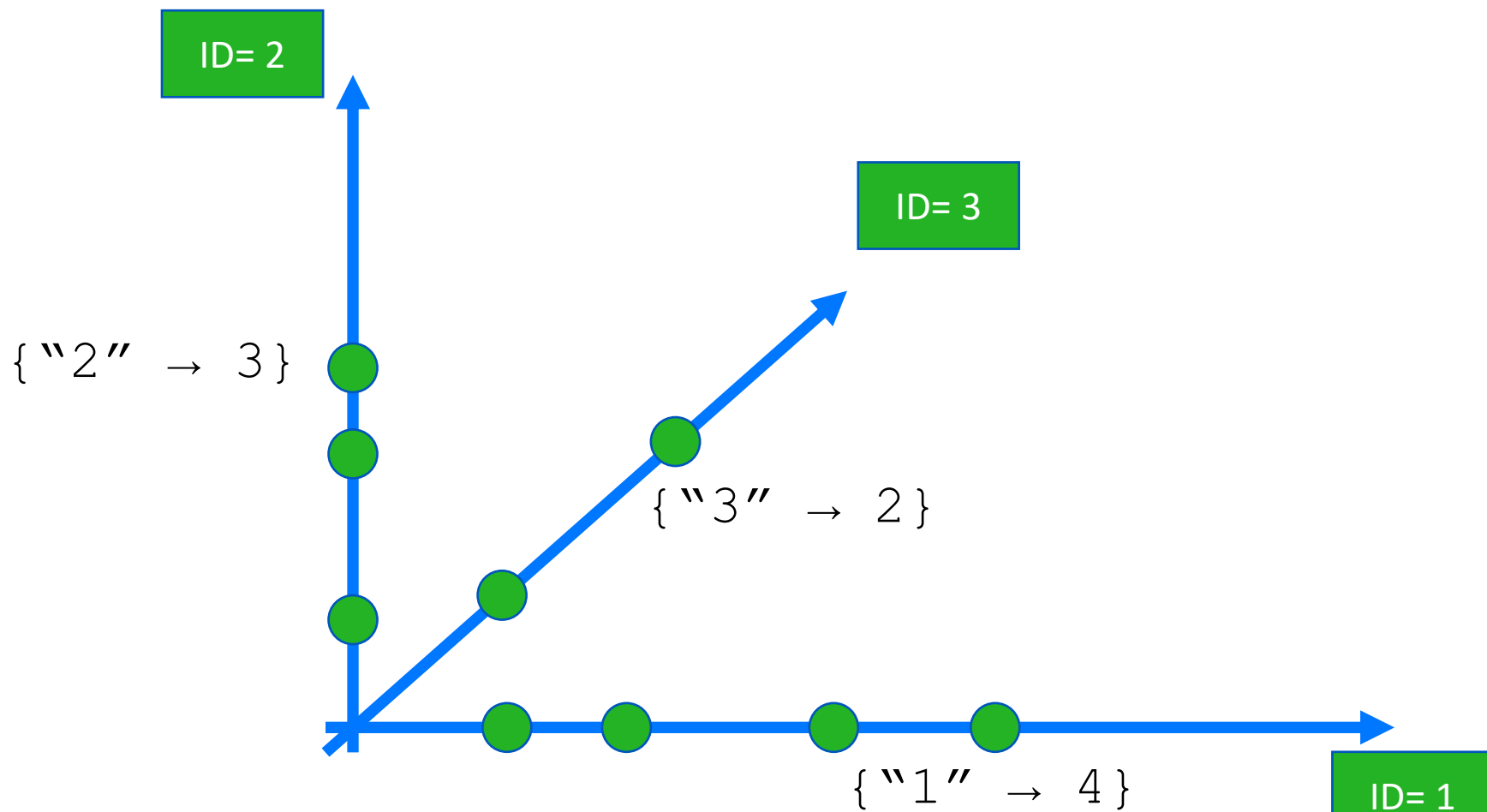
Векторные часы (vclock)



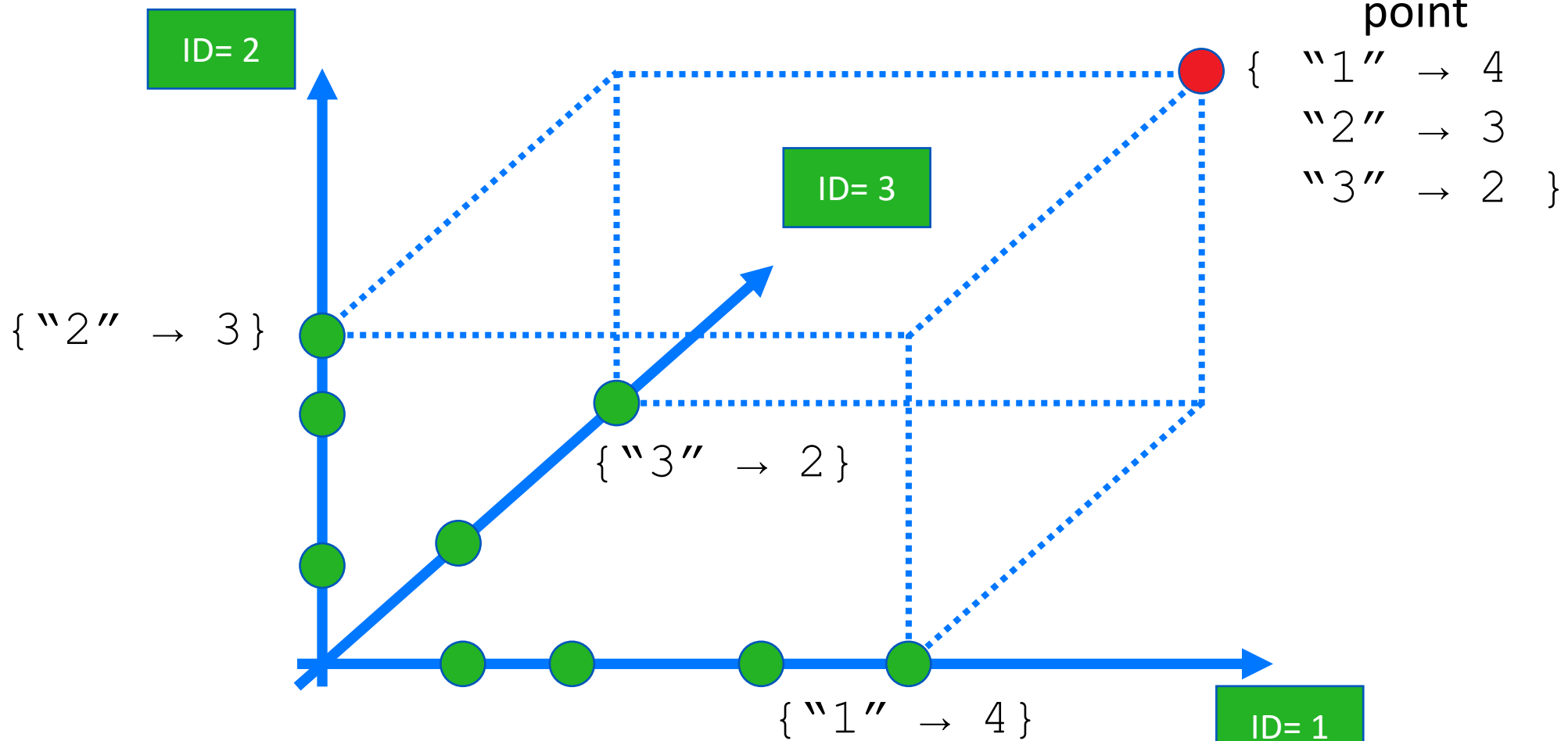
Векторные часы (vclock)



Векторные часы (vclock)



Векторные часы (vclock)



Репликация в Tarantool

- Основана на WAL
- Использует VClock для определения состояния
- Master-Master или Master-Replica
- Произвольная топология
- Синхронная и асинхронная*

* Может быть сконфигурирована per space

05

MVCC

Multi Version Concurrency Control
Борьба за чистоту чтений

Сохранность в памяти

Чтение «не ждёт»

Читаем из памяти

Запись «ждёт»

Пишем в память

Пишем на диск

Сохранность в памяти

Чтение «не ждёт»

Читаем из памяти

Читаем из памяти (“dirty”)

Запись «ждёт»

Пишем в память

Пишем на диск (долго)

Консистентность чтений

RAM

0

Disk

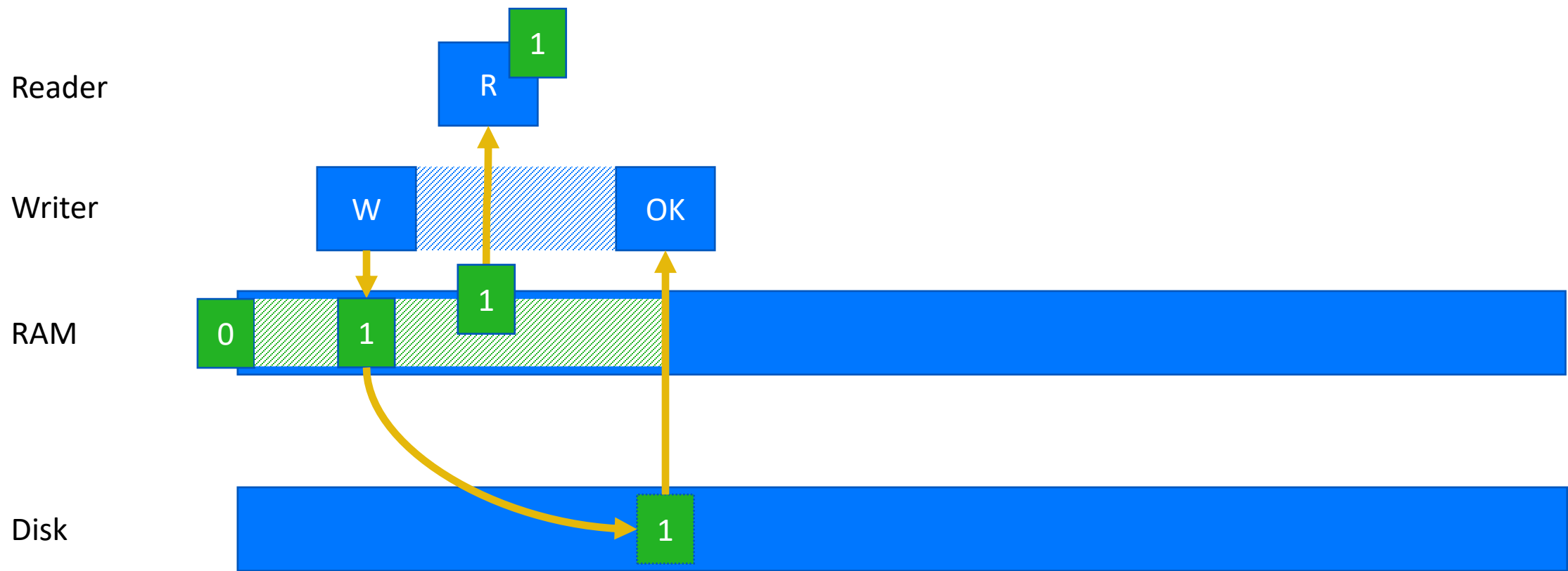
Запись происходит в память



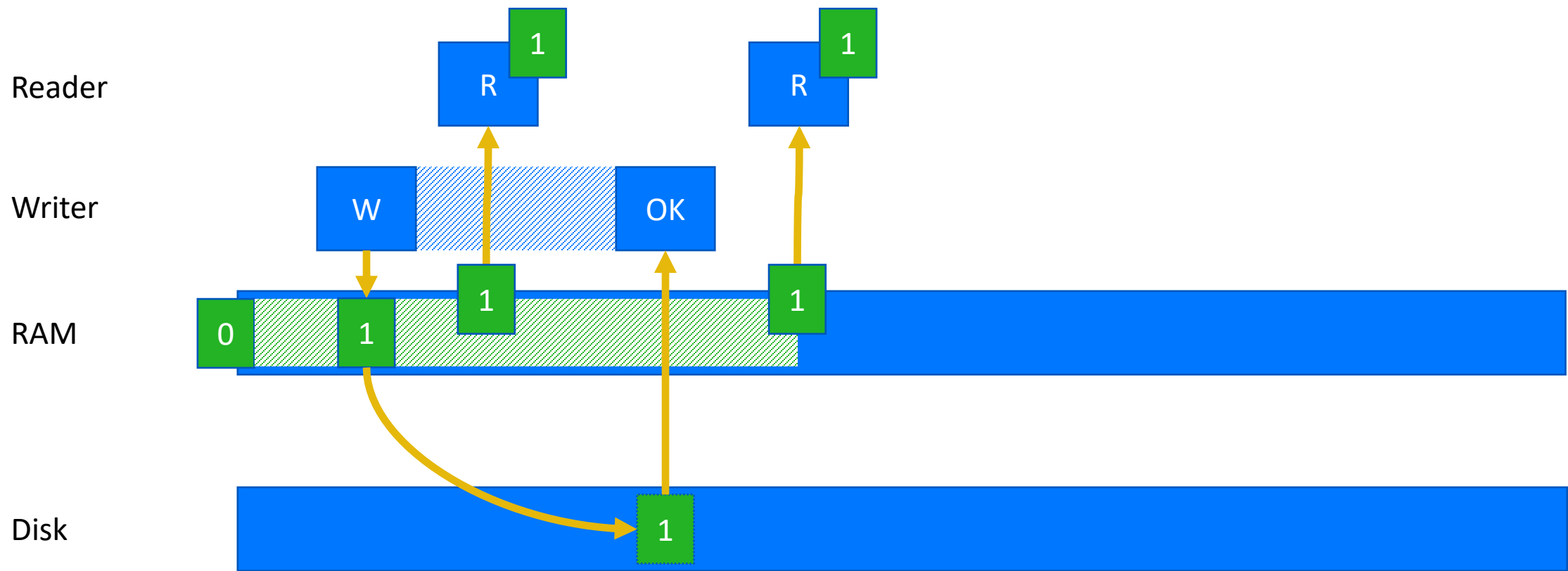
Чтение происходит из памяти



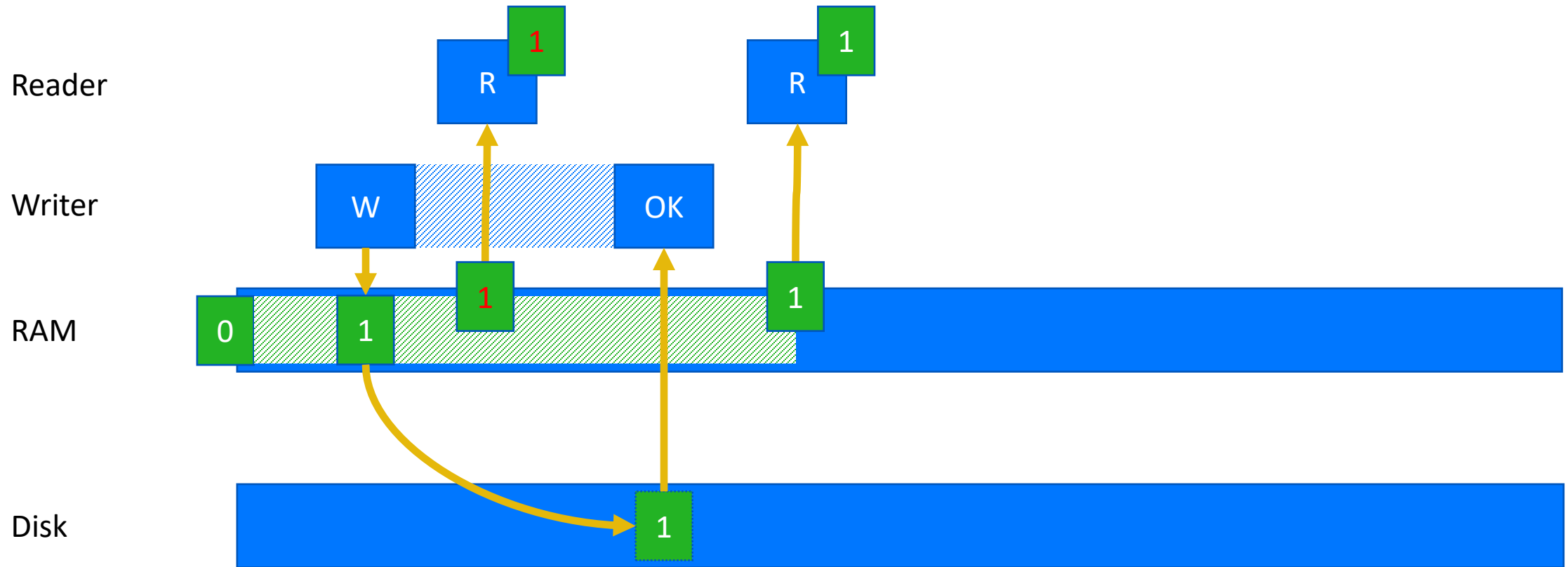
Запись на диск подтверждается



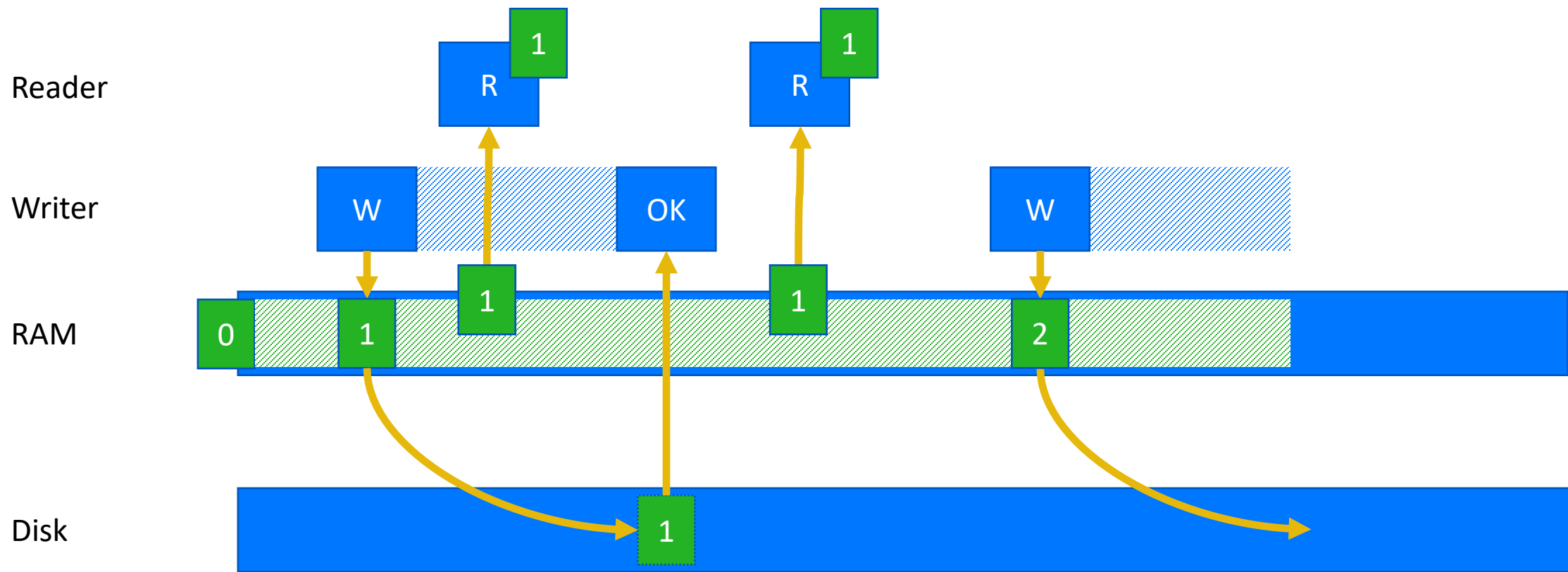
Запись на диск подтверждается



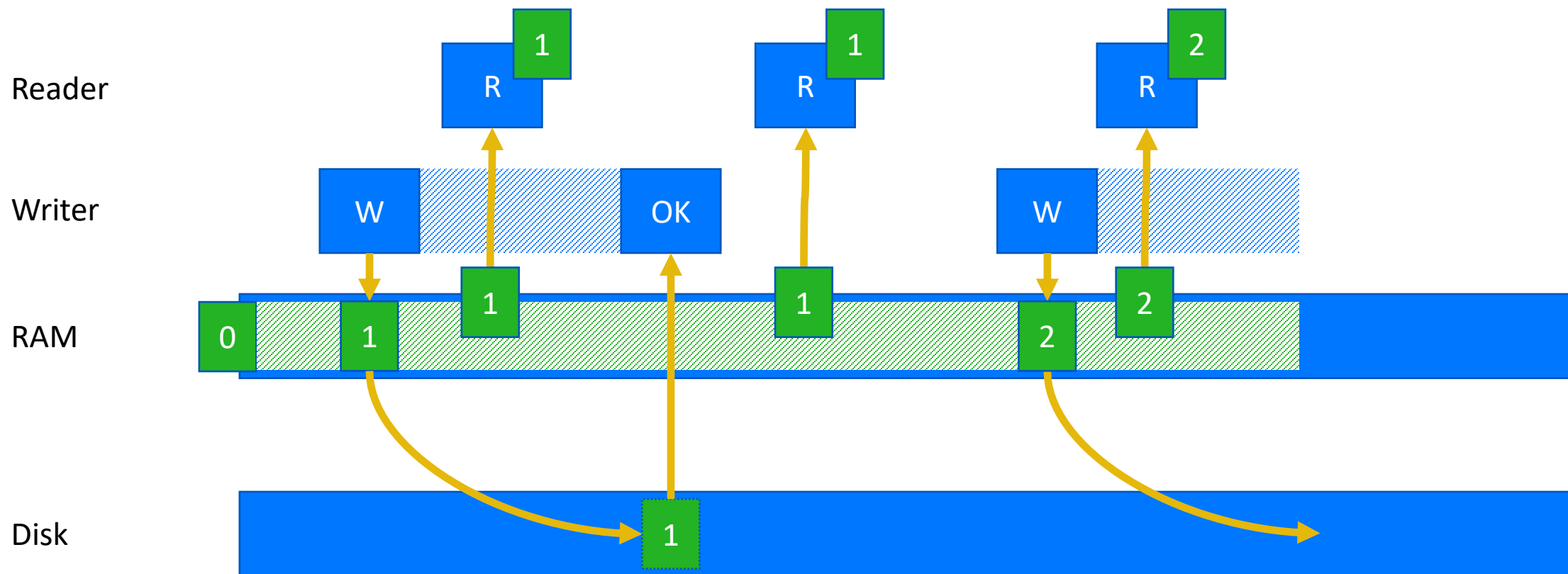
Некоторые чтения — до записи в WAL



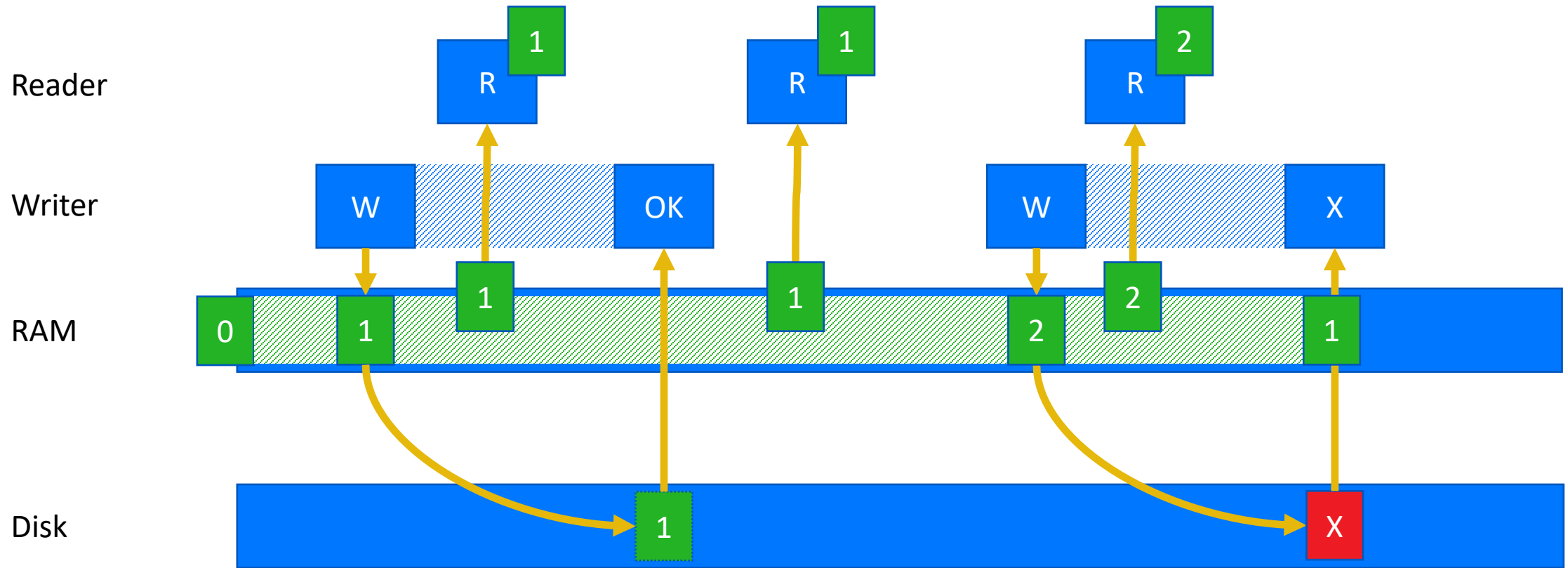
(не)консистентность чтений?



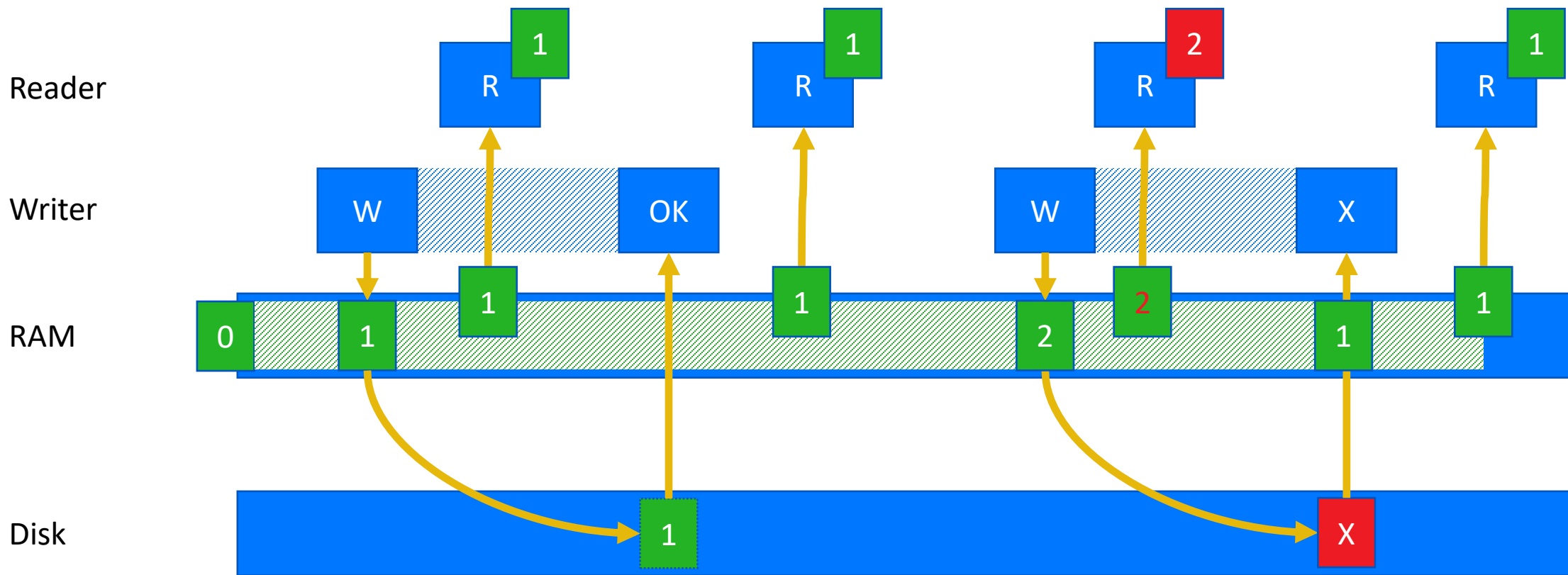
(не)консистентность чтений?



Ошибка диска откатывает транзакцию



Читатель видел «грязные» данные



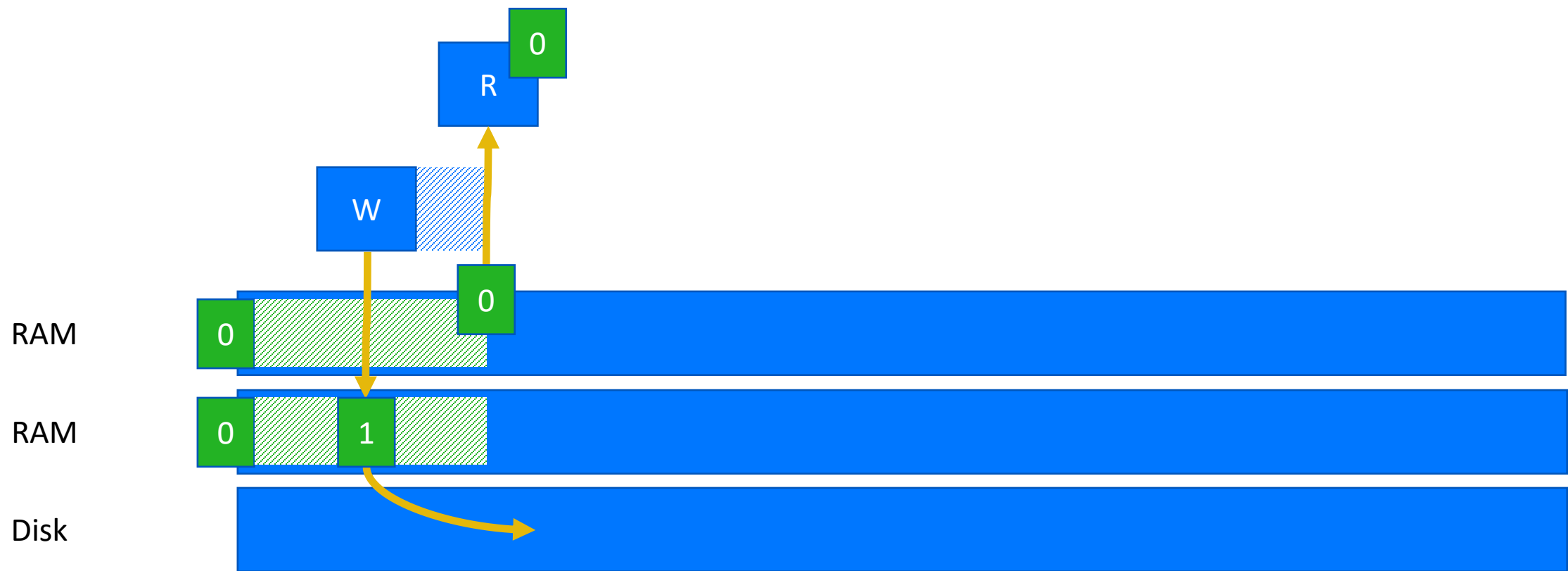
Разделить области чтения и записи



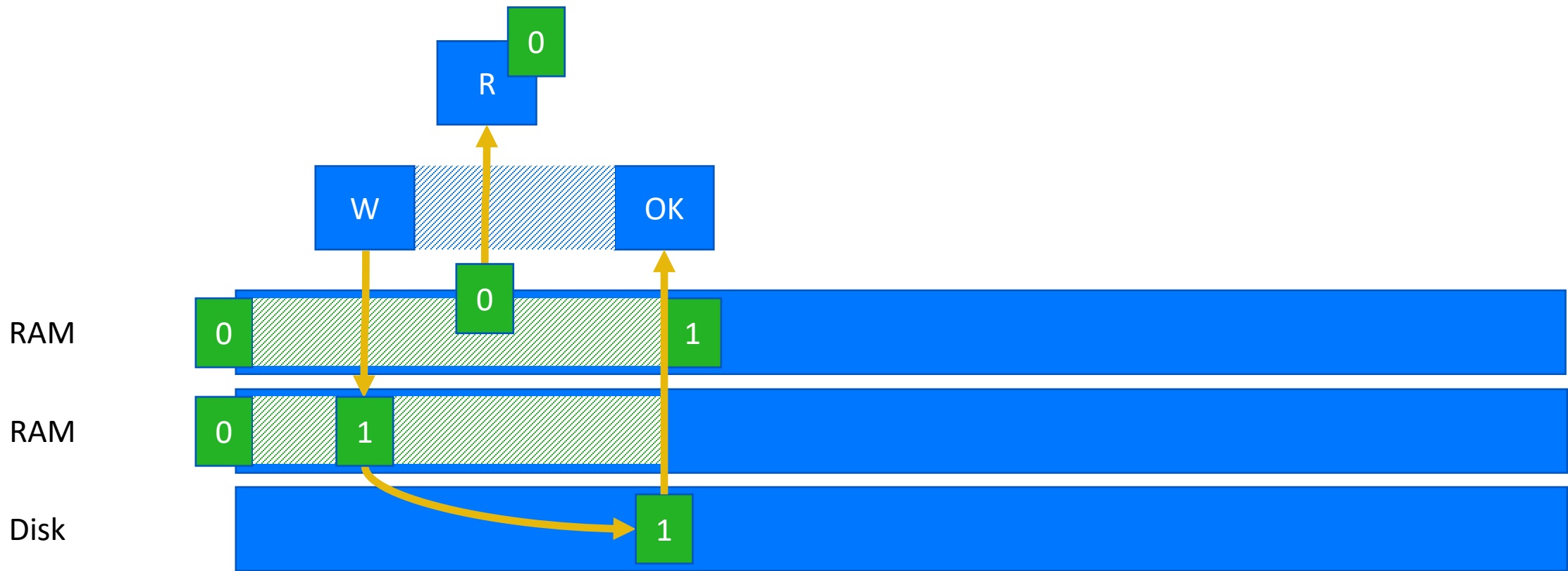
Разделить области чтения и записи



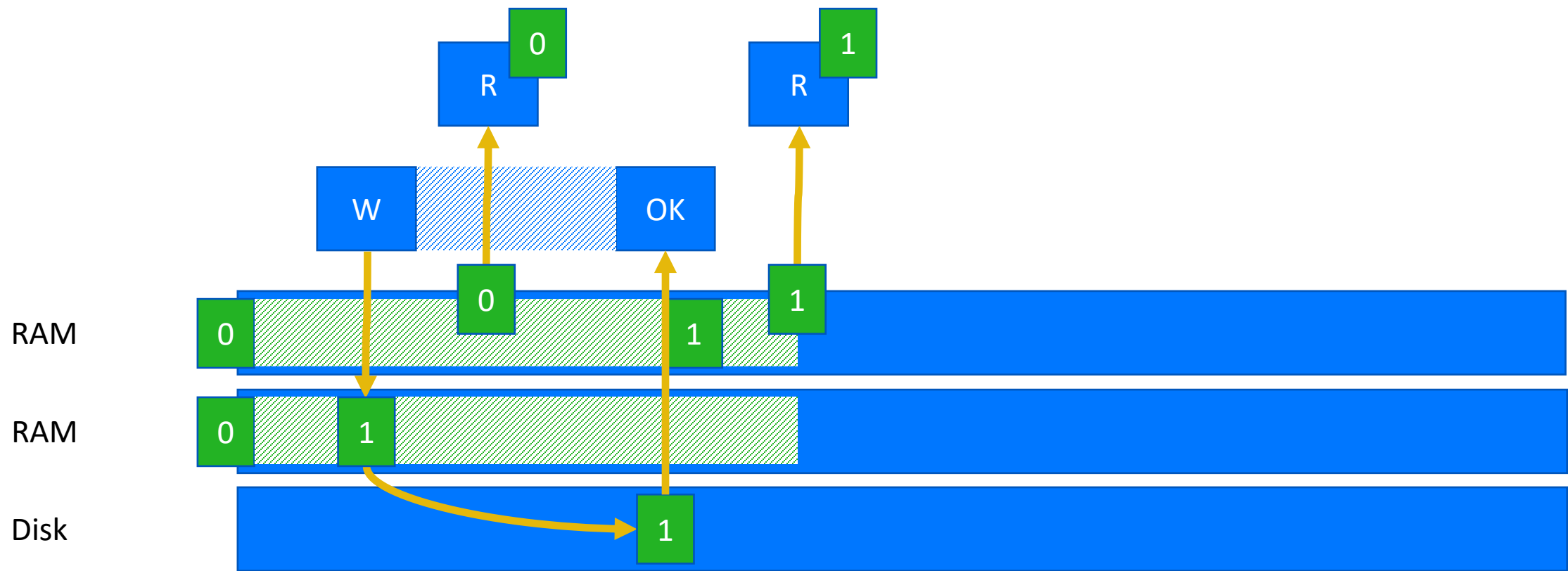
Чтение из read-view



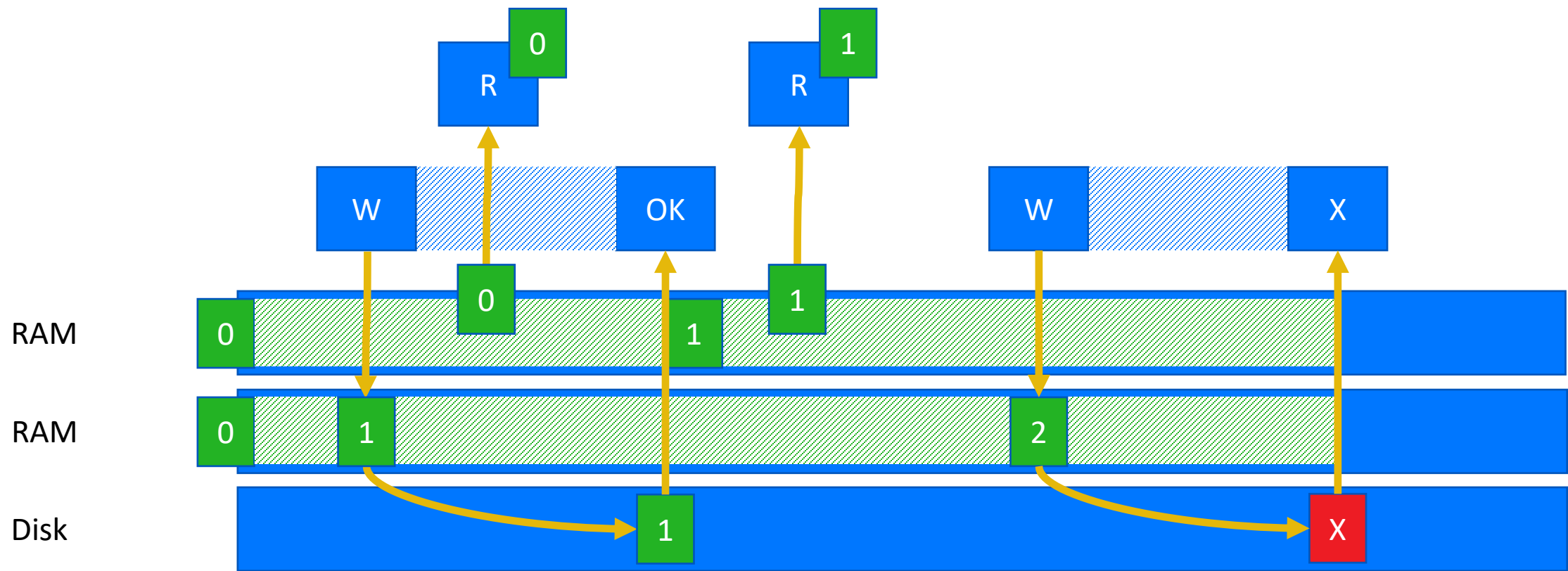
Чтение из read-view



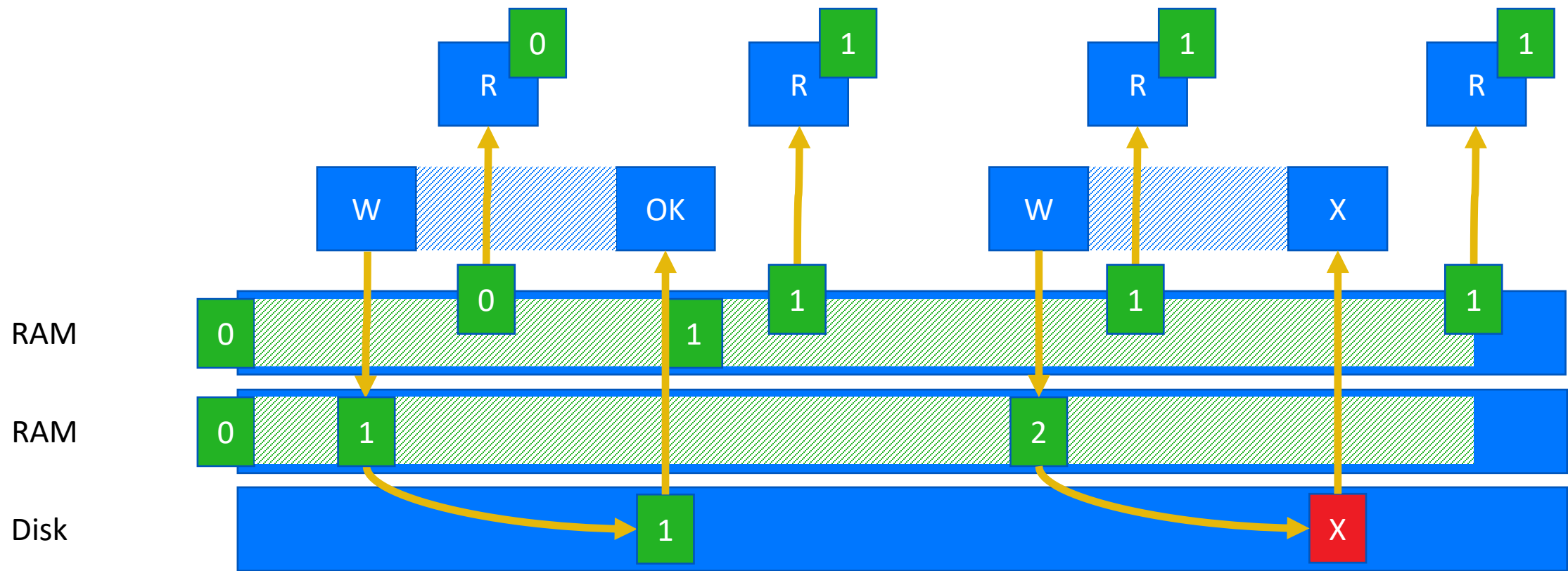
Чтение только подтверждённого



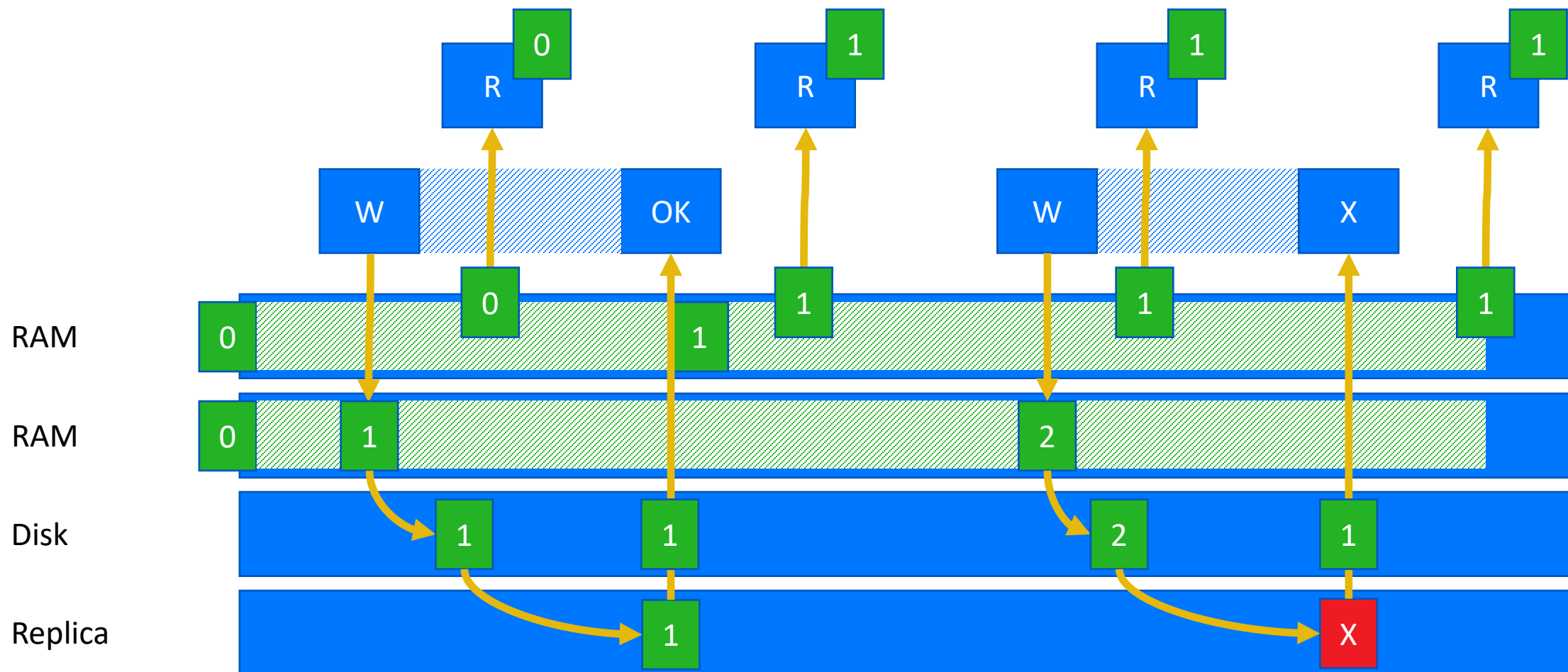
Чтение только подтверждённого



Чтение только подтверждённого



Чтение только отреплицированного



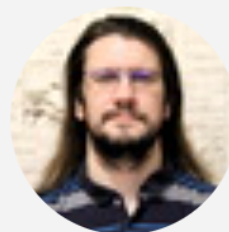
Подробнее об MVCC — завтра

h3: Яндекс трек

25 ноября
(завтра)

14:40 – 15:30

Как работает MVCC
в In-Memory-СУБД



Александр Ляпунов (Tarantool, VK)

Как устроен Tarantool?

- 100% данных находятся в памяти
- Доступ к данным из одного потока
- Доступ к данным только через индекс
- Изменения пишутся во Write Ahead Log (WAL)
- ACID обеспечивается однопоточностью и WAL
- Периодически сохраняются консистентный Snapshot
- WAL реплицируется
- Репликация гибридная (MM, MR, Sync, Async)

Спасибо за внимание!



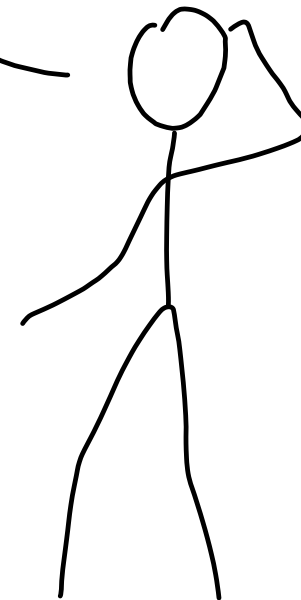
Оставьте
отзыв

Скачайте
презентацию



bx.vc/xxx

Вопросы?



Mons Anderson



166

Яндекс

